

facebook

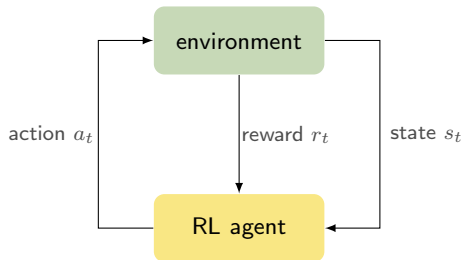
Artificial Intelligence Research

Exploration-Exploitation in Reinforcement Learning

Ronan Fruit^{*}, Alessandro Lazaric[†] and Matteo Pirodda[†]

Facebook AI Research[†] and INRIA Lille^{*}

Reinforcement Learning



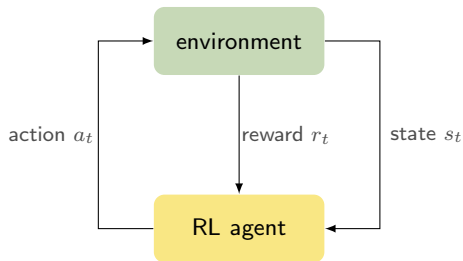
“**Reinforcement learning** is learning how to map states to actions so as to **maximize** a numerical **reward** signal in an unknown and **uncertain** environment.

In the most interesting and challenging cases, **actions** affect not only the immediate reward but also the **next situation** and all subsequent rewards (**delayed reward**).

The agent is not told which actions to take but it must discover which actions yield the most reward by trying them (**trial-and-error**).”

— Sutton and Barto [1998]

Reinforcement Learning



“**Reinforcement learning** is learning how to map states to actions so as to **maximize** a numerical **reward** signal in an unknown and **uncertain** environment.

In the most interesting and challenging cases, **actions** affect not only the immediate reward but also the **next situation** and all subsequent rewards (**delayed reward**).

The agent is not told which actions to take but it must discover which actions yield the most reward by trying them (**trial-and-error**).”

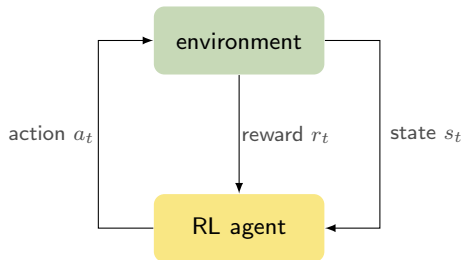
Exploration

— Sutton and Barto [1998]

Reinforcement Learning

Exploitation

“**Reinforcement learning** is learning how to map states to actions so as to maximize a numerical **reward** signal in an unknown and **uncertain** environment.



In the most interesting and challenging cases, **actions** affect not only the immediate reward but also the **next situation** and all subsequent rewards (**delayed reward**).

Exploration

The agent is not told which actions to take but it must discover which actions yield the most reward by trying them (**trial-and-error**).”

— Sutton and Barto [1998]

Disclaimer: the Real Title

Regret Minimization in Infinite-Horizon Finite Markov Decision Processes

Organization

- 1 Setting the Stage
- 2 Lower Bounds
- 3 Optimism in Face of Uncertainty
- 4 Posterior Sampling
- 5 Asymptotically Optimal Algorithms
- 6 Extensions and Other Settings
- 7 Conclusion

Website

<https://rlgammazero.github.io>

Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space \mathcal{S} , $|\mathcal{S}| = S < \infty$
- Action space \mathcal{A} , $|\mathcal{A}| = A < \infty$
- Transition distribution $p(\cdot | s, a) \in \Delta(\mathcal{S})$
- Reward distribution with expectation $r(s, a) \in [0, r_{\max}]$

Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space \mathcal{S} , $|\mathcal{S}| = S < \infty$
 - Action space \mathcal{A} , $|\mathcal{A}| = A < \infty$
- } finite
- Transition distribution $p(\cdot | s, a) \in \Delta(\mathcal{S})$
 - Reward distribution with expectation $r(s, a) \in [0, r_{\max}]$

Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space \mathcal{S} , $|\mathcal{S}| = S < \infty$
 - Action space \mathcal{A} , $|\mathcal{A}| = A < \infty$
- } **finite**
- Transition distribution $p(\cdot | s, a) \in \Delta(\mathcal{S})$ } **Markov**
 - Reward distribution with expectation $r(s, a) \in [0, r_{\max}]$

Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space \mathcal{S} , $|\mathcal{S}| = S < \infty$
 - Action space \mathcal{A} , $|\mathcal{A}| = A < \infty$
 - Transition distribution $p(\cdot | s, a) \in \Delta(\mathcal{S})$
 - Reward distribution with expectation $r(s, a) \in [0, r_{\max}]$
- } **finite**
- } **Markov**

👉 The process generates history $H_t = (s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t)$, with $s_{t+1} \sim p(\cdot | s_t, a_t)$

Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space \mathcal{S} , $|\mathcal{S}| = S < \infty$
 - Action space \mathcal{A} , $|\mathcal{A}| = A < \infty$
- } **finite**
- Transition distribution $p(\cdot | s, a) \in \Delta(\mathcal{S})$ } **Markov**
 - Reward distribution with expectation $r(s, a) \in [0, r_{\max}]$
- 👉 The process generates history $H_t = (s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t)$, with $s_{t+1} \sim p(\cdot | s_t, a_t)$

📖 In (contextual) bandit, actions do not influence the evolution of states

Policies

An agent acts according to a *policy*

	stationary	history-dependent
deterministic	$\pi : \mathcal{S} \rightarrow \mathcal{A}$	$\pi_t : \mathcal{H}_t \rightarrow \mathcal{A}$
stochastic	$\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$	$\pi_t : \mathcal{H}_t \rightarrow \Delta(\mathcal{A})$

Classification

An MDP M is

- *ergodic* if it is possible to go from any state to any other state under *any* deterministic stationary policy

$$\forall s, s', \forall \pi : \mathcal{S} \rightarrow \mathcal{A}, \exists t < \infty, \text{ s.t. } \mathbb{P}_{\pi}^M(s_t = s' | s_0 = s) > 0$$

- *communicating* if it is possible to go from any state to any other state under *a specific* deterministic stationary policy

$$\forall s, s', \exists \pi : \mathcal{S} \rightarrow \mathcal{A}, \exists t < \infty, \text{ s.t. } \mathbb{P}_{\pi}^M(s_t = s' | s_0 = s) > 0$$

👉 A communicating MDP has *finite diameter*

$$D_M = \max_{s, s' \in \mathcal{S}} \min_{\pi : \mathcal{S} \rightarrow \mathcal{A}} \mathbb{E}[T_{\pi}^M(s, s')]$$

Classification

An MDP M is

- *ergodic* if it is possible to go from any state to any other state under *any* deterministic stationary policy

$$\forall s, s', \forall \pi : \mathcal{S} \rightarrow \mathcal{A}, \exists t < \infty, \text{ s.t. } \mathbb{P}_{\pi}^M(s_t = s' | s_0 = s) > 0$$

- *communicating* if it is possible to go from any state to any other state under *a specific* deterministic stationary policy

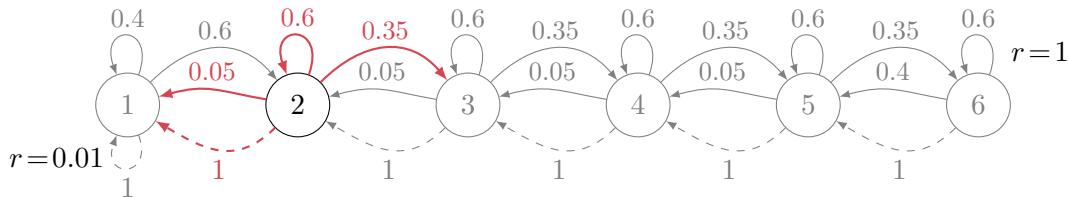
$$\forall s, s', \exists \pi : \mathcal{S} \rightarrow \mathcal{A}, \exists t < \infty, \text{ s.t. } \mathbb{P}_{\pi}^M(s_t = s' | s_0 = s) > 0$$

👉 A communicating MDP has *finite diameter*

$$D_M = \max_{s, s' \in \mathcal{S}} \underbrace{\min_{\pi : \mathcal{S} \rightarrow \mathcal{A}} \mathbb{E}[T_{\pi}^M(s, s')]}_{\text{shortest path}}$$

River Swim: Markov Decision Processes

Strehl and Littman [2008]



- $\mathcal{S} = \{1, 2, 3, 4, 5, 6\}$, $\mathcal{A} = \{L, R\}$
- $\pi_L(s) = L$, $\pi_R(s) = R$
- $M \oplus \pi_R$ is *ergodic* but $M \oplus \pi_L$ is *not ergodic*
- $T_{\pi_L}^M(6, 1) = 5$, $D_M = \mathbb{E}[T_{\pi_R}^M(1, 6)] \approx 14.7$

Gain and Bias

Gain of a deterministic stationary policy π

$$g_M^\pi(s) = \lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T r(s_t, a_t) \middle| s_0 = s, a_t = \pi(s_t) \right]$$

Bias of a deterministic stationary policy π

$$h_M^\pi(s) := C\text{-}\lim_{T \rightarrow \infty} \mathbb{E} \left[\sum_{t=1}^T (r(s_t, a_t) - g_M^\pi(s_t)) \middle| s_0 = s, a_t = \pi(s_t) \right]$$

Span of the bias function

$$\text{sp}(h_M^\pi) = \max_s h_M^\pi(s) - \min_s h_M^\pi(s)$$

Bellman operators

Bellman operator $L_M^a : \mathbb{R}^S \rightarrow \mathbb{R}^S$

$$= \sum_{s'} p(s'|s, a) h(s')$$

$$L_M^a h(s) = r(s, a) + p(\cdot|s, a)^\top h$$

Optimal Bellman operator $L_M^\star : \mathbb{R}^S \rightarrow \mathbb{R}^S$

$$L_M^\star h(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot|s, a)^\top h \right\}$$

Optimality gap of action a at s

$$\delta_M^\star(s, a) = L_M^\star h_M^\star(s) - L_M^a h_M^\star(s)$$

a.k.a. advantage function

Optimality

Optimal policy and *optimal gain*

$$\pi_M^* \in \arg \max_{\pi} g_M^{\pi}(s) \quad g_M^* = g_M^{\pi^*}(s) \quad \forall s \in \mathcal{S}$$

Optimality equation

$$h_M^*(s) + g_M^* = L_M^* h_M^*(s)$$

Greedy policy w.r.t. h_M^* is optimal

$$\pi_M^*(s) \in \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^{\top} h_M^* \right\}$$

Set of optimal actions in state s

$$\Pi_M^*(s) = \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^{\top} h_M^* \right\}$$

Optimality

deterministic stationary

Optimal policy and *optimal gain*

$$\pi_M^* \in \arg \max_{\pi} g_M^{\pi}(s) \quad g_M^* = g_M^{\pi^*}(s) \quad \forall s \in \mathcal{S}$$

Optimality equation

$$h_M^*(s) + g_M^* = L_M^* h_M^*(s)$$

Greedy policy w.r.t. h_M^* is optimal

$$\pi_M^*(s) \in \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^{\top} h_M^* \right\}$$

Set of optimal actions in state s

$$\Pi_M^*(s) = \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^{\top} h_M^* \right\}$$

Optimality

deterministic stationary

Optimal policy and *optimal gain*

constant gain*

$$\pi_M^* \in \arg \max_{\pi} g_M^{\pi}(s) \quad g_M^* = g_M^{\pi^*}(s) \quad \forall s \in \mathcal{S}$$

Optimality equation

$$h_M^*(s) + g_M^* = L_M^* h_M^*(s)$$

Greedy policy w.r.t. h_M^* is optimal

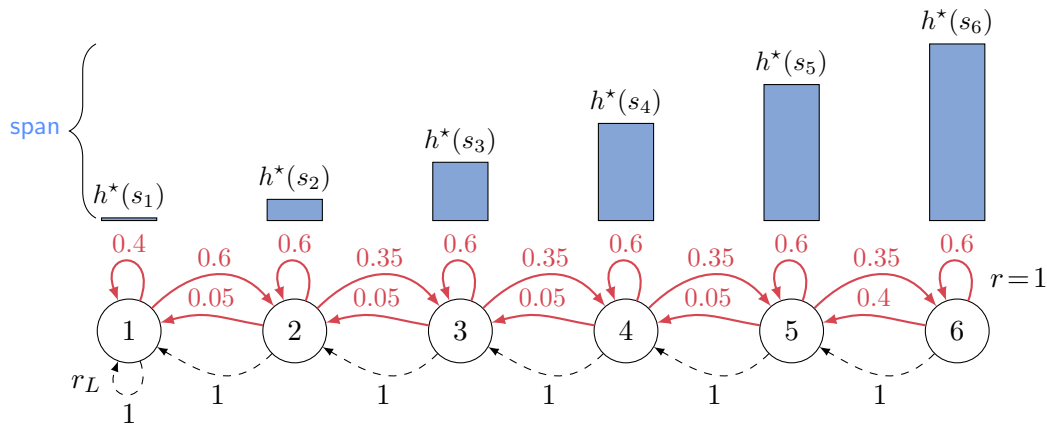
$$\pi_M^*(s) \in \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^{\top} h_M^* \right\}$$

Set of optimal actions in state s

$$\Pi_M^*(s) = \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^{\top} h_M^* \right\}$$

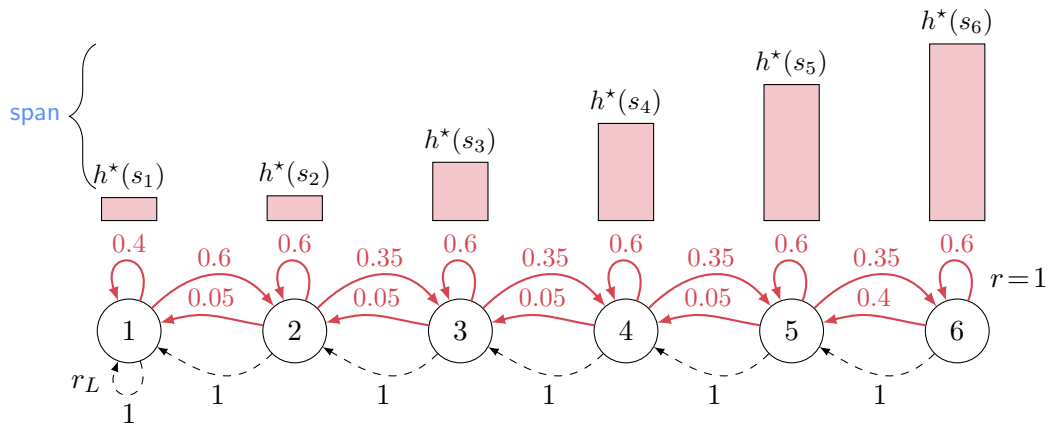
*In communicating MDPs

River Swim: Optimality



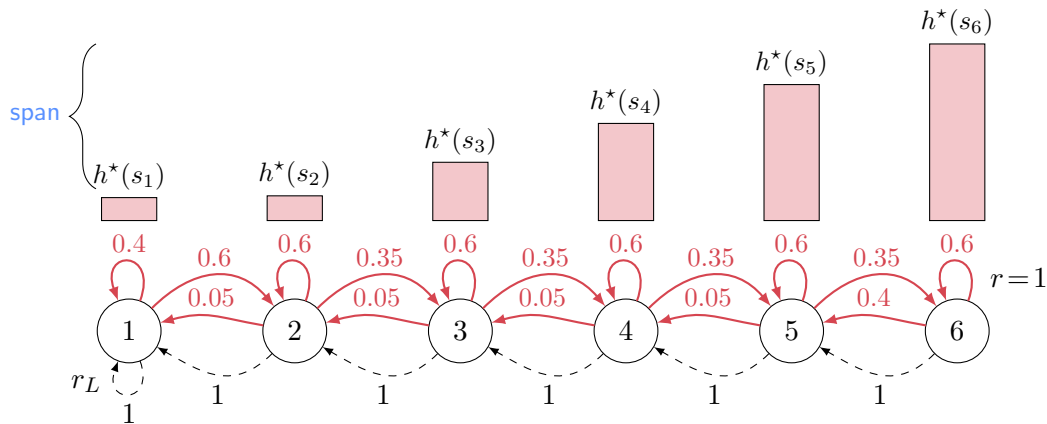
- $\pi^* = \pi_R$
- If $r_L = 0.01$, $g^* \approx 0.43$, $\text{sp}(h^*) \approx 6.4$

River Swim: Optimality



- $\pi^* = \pi_R$
- If $r_L = 0.01$, $g^* \approx 0.43$, $\text{sp}(h^*) \approx 6.4$
- If $r_L = 0.4$, $g^* \approx 0.43$, $\text{sp}(h^*) \approx 5.5$

River Swim: Optimality



- $\pi^* = \pi_R$
 - If $r_L = 0.01$, $g^* \approx 0.43$, $\text{sp}(h^*) \approx 6.4$
 - If $r_L = 0.4$, $g^* \approx 0.43$, $\text{sp}(h^*) \approx 5.5$
- } D is constant

Value Iteration

initialize $v_0(s) = 0 \quad \forall s \in \mathcal{S}, n = 0, \varepsilon$

repeat

for $s \in \mathcal{S}$ **do**

$v_{n+1}(s) = L_M^* v_n(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^\top v_n \right\}$

end

$n = n + 1$

until $sp(v_{n+1} - v_n) < \varepsilon$

return greedy policy

$$\pi_\varepsilon(s) = \arg \max_{a \in \mathcal{A}} L_M^a v_n(s) = \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + p(\cdot | s, a)^\top v_n \right\}$$

Value Iteration

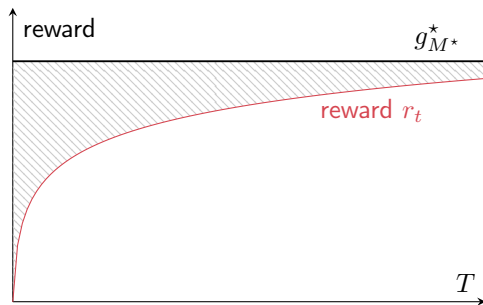
Theorem (Thm. 8.5.5 [Puterman, 1994])

In any communicating MDP M , value iteration is such that

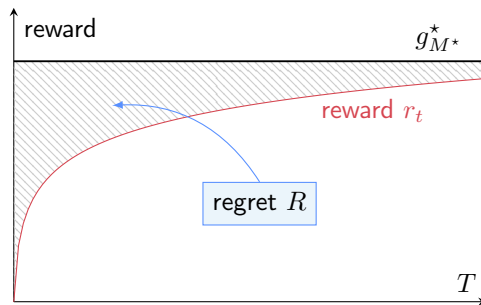
- *convergence: for any ε , there exists n_ε s.t. the stopping condition is met*
- *optimality: policy π_ε is ε -optimal*

$$g_M^{\pi_\varepsilon}(s) \geq g_M^* - \varepsilon$$

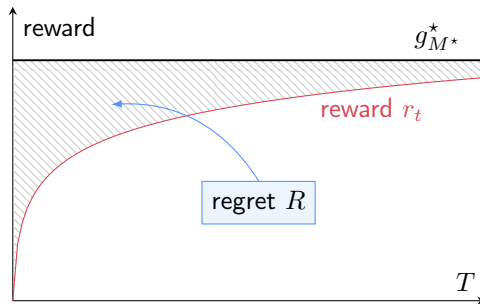
Regret Minimization



Regret Minimization



Regret Minimization



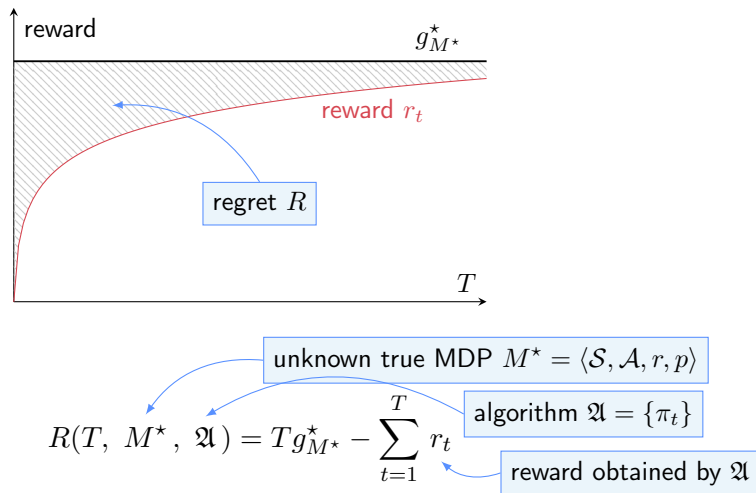
$$R(T, M^*, \mathfrak{A}) = T g_{M^*}^* - \sum_{t=1}^T r_t$$

unknown true MDP $M^* = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

algorithm $\mathfrak{A} = \{\pi_t\}$

reward obtained by \mathfrak{A}

Regret Minimization



Expected regret w.r.t. randomness of s_t , r_t , and (possibly) \mathfrak{A}

$$\bar{R}(T, M^*, \mathfrak{A}) = \mathbb{E}[R(T, M^*, \mathfrak{A})]$$

- 1 Setting the Stage
- 2 Lower Bounds
- 3 Optimism in Face of Uncertainty
- 4 Posterior Sampling
- 5 Asymptotically Optimal Algorithms
- 6 Extensions and Other Settings
- 7 Conclusion

Problem-Dependent Lower Bound

Let $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$ and $M' = \langle \mathcal{S}, \mathcal{A}, r, p' \rangle$

- *Difference* between M and M' at s, a (w.l.o.g. assuming reward known)

$$\text{KL}_{M,M'}(s, a) = \text{KL}(p(\cdot|s, a) \| p'(\cdot|s, a))$$

- *Set of alternative* (confusing) models w.r.t. M

same everywhere but in (s, a)

$$\mathcal{M}_M^{\text{alt}}(s, a) = \left\{ M' : p'(\cdot|s', a') = p(\cdot|s', a'), \text{ for all } (s', a') \neq (s, a), \right. \\ \left. a \notin \Pi_M^*(s), a \in \Pi_{M'}^*(s) \right\}$$

sub-optimal in M

optimal in M'

Problem-Dependent Lower Bound

Theorem (Thm. 1 Burnetas and Katehakis [1997], Thm. 2 Ok et al. [2018])

Let \mathfrak{A} be s.t. $\bar{R}(T, M, \mathfrak{A}) = o(T^\alpha)$ for all $\alpha > 0$ and *ergodic* MDP M . For any *ergodic* MDP M^* with $r_{\max} = 1$, the expected regret is lower bounded as

$$\liminf_{T \rightarrow \infty} \frac{\bar{R}(T, M^*, \mathfrak{A})}{\log T} \geq K_{M^*}$$

where

$$K_{M^*} = \inf_{\eta \geq 0} \sum_{s,a} \eta(s,a) \delta_{M^*}^*(s,a)$$

$$\text{s.t. } \sum_{s,a} \eta(s,a) \text{KL}_{M^*,M}(s,a) \geq 1 \quad \forall M \in \mathcal{M}_{M^*}^{\text{alt}}(s,a)$$

cumulative regret

"evidence" of difference between M^* and M

Problem-Dependent Lower Bound

Theorem (Thm. 1 Burnetas and Katehakis [1997], Thm. 2 Ok et al. [2018])

Let \mathfrak{A} be s.t. $\bar{R}(T, M, \mathfrak{A}) = o(T^\alpha)$ for all $\alpha > 0$ and *ergodic* MDP M . For any *ergodic* MDP M^* with $r_{\max} = 1$, the expected regret is lower bounded as

$$\liminf_{T \rightarrow \infty} \frac{\bar{R}(T, M^*, \mathfrak{A})}{\log T} \geq K_{M^*}$$

where

$$K_{M^*} = \inf_{\eta \geq 0} \sum_{s,a} \eta(s,a) \delta_{M^*}^*(s,a)$$

$$\text{s.t. } \sum_{s,a} \eta(s,a) \text{KL}_{M^*,M}(s,a) \geq 1 \quad \forall M \in \mathcal{M}_{M^*}^{\text{alt}}(s,a)$$

cumulative regret

“evidence” of difference between M^* and M

Similar to [Lai and Robbins, 1985] for MAB but alternative models and regret are different.

Problem-Dependent Lower Bound

Theorem (Thm. 1 Burnetas and Katehakis [1997], Thm. 2 Ok et al. [2018])

Let \mathfrak{A} be s.t. $\bar{R}(T, M, \mathfrak{A}) = o(T^\alpha)$ for all $\alpha > 0$ and *ergodic* MDP M . For any *ergodic* MDP M^* with $r_{\max} = 1$, the expected regret is lower bounded as

$$\liminf_{T \rightarrow \infty} \frac{\bar{R}(T, M^*, \mathfrak{A})}{\log T} \geq K_{M^*}$$

where

$$K_{M^*} \leq 2 \frac{(C + 1)^2}{\min_{s,a} \delta_{M^*}(s, a)} SA \quad C = sp(h_{M^*}^*)$$

Minimax Lower Bound

Theorem (Thm. 5 Jaksch et al. [2010])

For any *communicating* MDP M^* with $r_{\max} = 1$, $S, A \geq 10$, $D \geq 20 \log_A S$, any algorithm \mathfrak{A} at any time $T \geq DSA$ suffers a regret

$$\sup_{M^*} \bar{R}(T, M^*, \mathfrak{A}) \geq 0.015 \sqrt{DSAT}$$

Minimax Lower Bound

Theorem (Thm. 5 Jaksch et al. [2010])

For any *communicating* MDP M^\star with $r_{\max} = 1$, $S, A \geq 10$, $D \geq 20 \log_A S$, any algorithm \mathfrak{A} at any time $T \geq DSA$ suffers a regret

$$\sup_{M^\star} \bar{R}(T, M^\star, \mathfrak{A}) \geq 0.015 \sqrt{DSAT}$$

 In MAB $\Omega(\sqrt{AT})$ since $D = 1$ and $S = 1$.

Open Questions

C could be arbitrarily large
($C = \infty$ for non ergodic)

- 1 *Asymptotic* regime and *ergodicity* assumption

$$\mathbb{P}_M^\pi[N_T(s) \geq \rho T] \geq 1 - C \exp(-\rho T/2) \quad [\text{Prop.2 Burnetas and Katehakis [1997]]]$$

- 2 *Span vs. diameter*

$D = 2\text{sp}(h^*)$ in the proof

$$\overline{R}(T, M^*, \mathfrak{A}) \geq 0.015 \sqrt{D \text{ SAT}}$$

- 3 *Number of states vs branching factor* $\Gamma = \max_{s,a} |\text{supp}(p(\cdot|s, a))|$

$$\overline{R}(T, M^*, \mathfrak{A}) \geq 0.015 \sqrt{D S A T}$$

$\Gamma = 2$ in the proof

- 1 Setting the Stage
- 2 Lower Bounds
- 3 Optimism in Face of Uncertainty
- 4 Posterior Sampling
- 5 Asymptotically Optimal Algorithms
- 6 Extensions and Other Settings
- 7 Conclusion

The Optimism Principle: Intuition



OPTIMISM
It's the best way to see life.

The Optimism Principle: Intuition

Exploration vs. Exploitation

The Optimism Principle: Intuition

Exploration vs. Exploitation

Optimism in Face of Uncertainty

When you are uncertain, consider the **best possible world** (reward-wise)

The Optimism Principle: Intuition

Exploration vs. Exploitation

Optimism in Face of Uncertainty

When you are uncertain, consider the **best possible world** (reward-wise)

If the best possible world is **correct**

⇒ **no regret**

Exploitation

If the best possible world is **wrong**

⇒ **learn useful information**

Exploration

The Optimism Principle: Intuition

Exploration vs. Exploitation

Optimism in gain

Optimism in Face of Uncertainty

When you are uncertain, consider the **best possible world** (reward-wise)

If the best possible world is **correct**

⇒ **no regret**

Exploitation

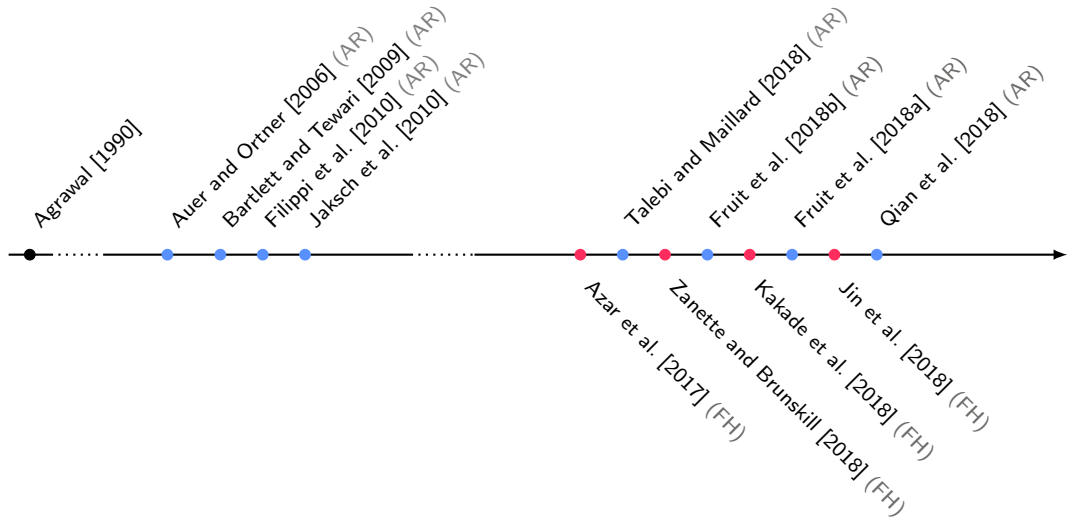
If the best possible world is **wrong**

⇒ **learn useful information**

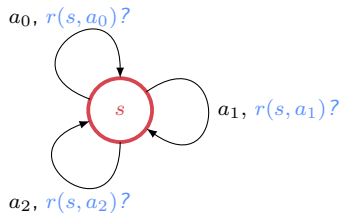
Exploration

History: OFU for Regret Minimization in RL

FH: finite-horizon
AR: average reward



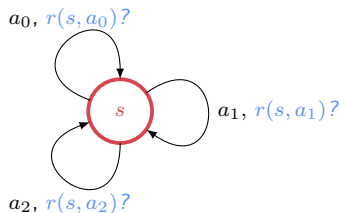
Gain Optimism: Example



■ Deterministic *policies*:

- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

Gain Optimism: Example



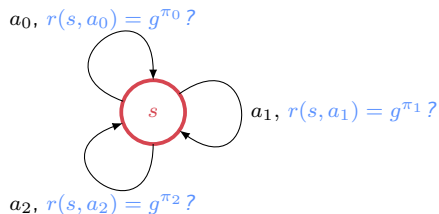
■ Deterministic *policies*:

- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

■ Optimism

$$\tilde{\pi} = \arg \max_{\pi_i} \text{UCB}(g^{\pi_i})$$

Gain Optimism: Example



■ Deterministic *policies*:

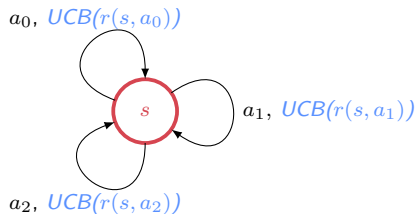
- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

■ Reward $r(s, a_i) = \text{gain } g^{\pi_i}$

■ Optimism

$$\tilde{\pi} = \arg \max_{\pi_i} \text{UCB}(g^{\pi_i})$$

Gain Optimism: Example



■ Deterministic *policies*:

- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

■ Reward $r(s, a_i) = \text{gain } g^{\pi_i}$

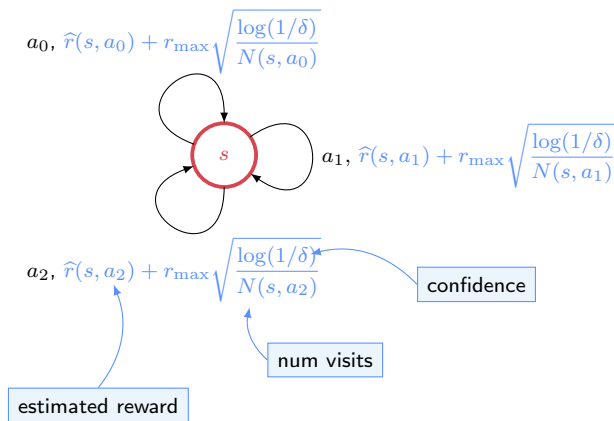
■ Upper confidence bound

$$UCB(g^{\pi_i}) = UCB(r(s, a_i))$$

■ Optimism

$$\tilde{\pi} = \arg \max_{\pi_i} UCB(g^{\pi_i})$$

Gain Optimism: Example



■ Deterministic *policies*:

- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

■ Reward $r(s, a_i) = \text{gain } g^{\pi_i}$

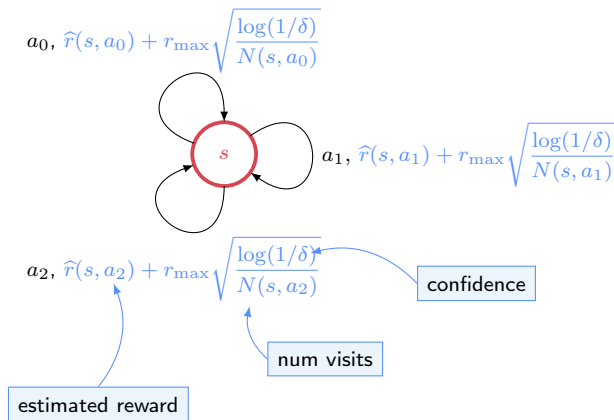
■ Upper confidence bound

$$\text{UCB}(g^{\pi_i}) = \text{UCB}(r(s, a_i))$$

■ Optimism

$$\tilde{\pi} = \arg \max_{\pi_i} \text{UCB}(g^{\pi_i})$$

Gain Optimism: Example



■ Deterministic *policies*:

- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

■ Reward $r(s, a_i) = \text{gain } g^{\pi_i}$

■ Upper confidence bound

$$\text{UCB}(g^{\pi_i}) = \text{UCB}(r(s, a_i))$$

■ Optimism

$$\tilde{\pi} = \arg \max_{\pi_i} \text{UCB}(g^{\pi_i})$$

👍 UCB algorithm (Bandit)

Gain Optimism: Implementation

Tentative algorithm

Observe s_1

for $t = 1, 2, \dots$ **do**

Compute $\pi_t \leftarrow \arg \max_{\pi} UCB_t(g^{\pi})$

 Take action $a_t = \pi_t(s_t)$

 Observe reward r_t and next state s_{t+1}

 Compute $UCB_{t+1}(g^{\pi})$ for all π based on $UCB_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

Gain Optimism: Implementation

Tentative algorithm

Observe s_1

for $t = 1, 2, \dots$ **do**

 Compute $\pi_t \leftarrow \arg \max_{\pi} UCB_t(g^{\pi})$

 Take action $a_t = \pi_t(s_t)$

 Observe reward r_t and next state s_{t+1}

 Compute $UCB_{t+1}(g^{\pi})$ for all π based on $UCB_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

 *3 major issues:*

- *Upper confidence bounds*: construct $UCB_t(g^{\pi})$ with unknown dynamics
- *Computational complexity*: exponential number of policies
- *Frequent policy update*: inefficient exploration

Gain Optimism: Implementation

Tentative algorithm

Observe s_1

for $t = 1, 2, \dots$ do

 Compute $\pi_t \leftarrow \arg \max_{\pi} UCB_t(g^{\pi})$

 Take action $a_t = \pi_t(s_t)$

 Observe reward r_t and next state s_{t+1}

 Compute $UCB_{t+1}(g^{\pi})$ for all π based on $UCB_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

 *3 major issues:*

- *Upper confidence bounds*: construct $UCB_t(g^{\pi})$ with unknown dynamics
- *Computational complexity*: exponential number of policies
- *Frequent policy update*: inefficient exploration

Bounded Parameter MDP: Definition

Bounded parameter MDP [Strehl and Littman, 2008]

$$\mathcal{M}_t = \left\{ \langle \mathcal{S}, \mathcal{A}, r, p \rangle : r(s, a) \in B_t^r(s, a), p(\cdot | s, a) \in B_t^p(s, a), \forall (s, a) \in \mathcal{S} \times \mathcal{A} \right\}$$

Compact *confidence sets*

$$B_t^r(s, a) := \left[\hat{r}_t(s, a) - \beta_t^r(s, a), \hat{r}_t(s, a) + \beta_t^r(s, a) \right]$$

$$B_t^p(s, a) := \left\{ p(\cdot | s, a) \in \Delta(\mathcal{S}) : \|p(\cdot | s, a) - \hat{p}_t(\cdot | s, a)\|_1 \leq \beta_t^p(s, a) \right\}$$

Bounded Parameter MDP: Definition

Bounded parameter MDP [Strehl and Littman, 2008]

$$\mathcal{M}_t = \left\{ \langle \mathcal{S}, \mathcal{A}, r, p \rangle : r(s, a) \in B_t^r(s, a), p(\cdot | s, a) \in B_t^p(s, a), \forall (s, a) \in \mathcal{S} \times \mathcal{A} \right\}$$

Compact *confidence sets*

$$B_t^r(s, a) := \left[\hat{r}_t(s, a) - \beta_t^r(s, a), \hat{r}_t(s, a) + \beta_t^r(s, a) \right]$$

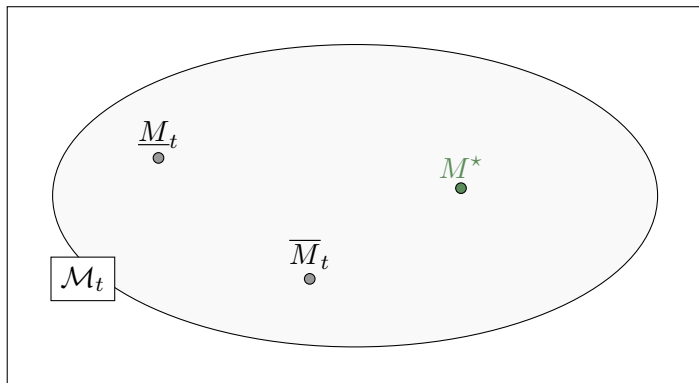
$$B_t^p(s, a) := \left\{ p(\cdot | s, a) \in \Delta(\mathcal{S}) : \|p(\cdot | s, a) - \hat{p}_t(\cdot | s, a)\|_1 \leq \beta_t^p(s, a) \right\}$$

Confidence bounds based on [Hoeffding, 1963] and [Weissman et al., 2003]

$$\beta_t^r(s, a) \propto \sqrt{\frac{\log(N_t(s, a)/\delta)}{N_t(s, a)}}$$

$$\beta_t^p(s, a) \propto \sqrt{\frac{S \log(N_t(s, a)/\delta)}{N_t(s, a)}}$$

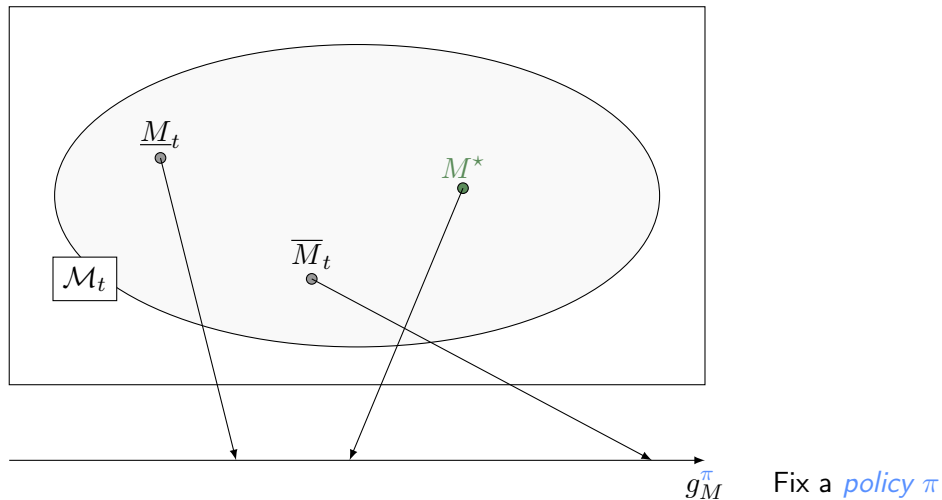
Bounded Parameter MDP: Optimism



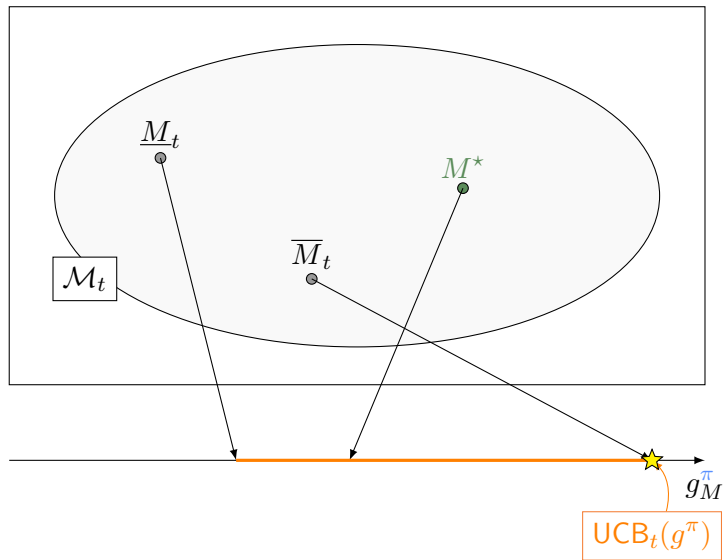
→ g_M^π

Fix a *policy* π

Bounded Parameter MDP: Optimism

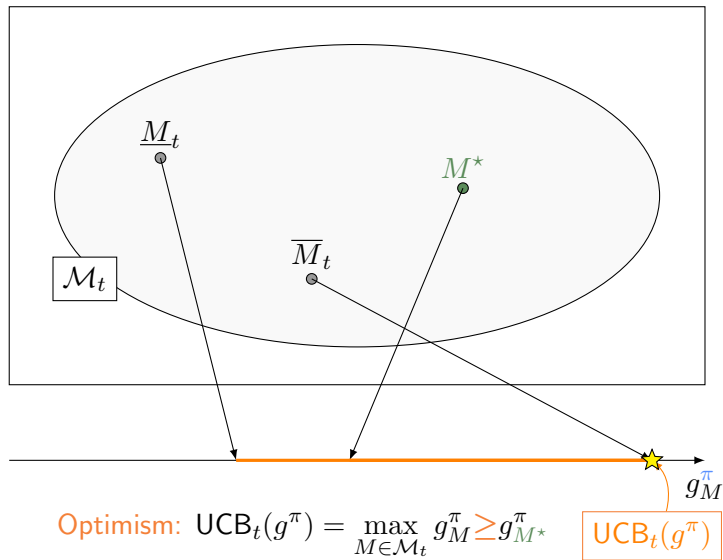


Bounded Parameter MDP: Optimism



Fix a *policy* π

Bounded Parameter MDP: Optimism



Fix a *policy* π

Gain Optimism: Implementation

Tentative algorithm

Observe state s_1

for $t = 1, 2, \dots$ do

 Compute $\pi_t \leftarrow \arg \max_{\pi} UCB_t(g^{\pi})$

 Take action $a_t = \pi_t(s_t)$

 Observe reward r_t and next state s_{t+1}

 Compute $UCB_{t+1}(g^{\pi})$ for all π based on $UCB_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

 *3 major issues:*

- *Upper confidence bounds*: construct $UCB_t(g^{\pi})$ with unknown dynamics? ✓
- *Computational complexity*: exponential number of policies
- *Frequent policy update*: inefficient exploration

Gain Optimism: Implementation

Tentative algorithm

Observe state s_1

for $t = 1, 2, \dots$ **do**

Compute $\pi_t \leftarrow \arg \max_{\pi} UCB_t(g^{\pi})$





Take action $a_t = \pi_t(s_t)$

Observe reward r_t and next state s_{t+1}

Compute $UCB_{t+1}(g^{\pi})$ for all π based on $UCB_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

 *3 major issues:*

-  *Upper confidence bounds:* construct $UCB_t(g^{\pi})$ with unknown dynamics? 
-  *Computational complexity:* exponential number of policies
-  *Frequent policy update:* inefficient exploration

Gain Optimism: Implementation

Tentative algorithm

Observe state s_1

for $t = 1, 2, \dots$ **do**

 Compute $\pi_t \leftarrow \arg \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\}$

 Take action $a_t = \pi_t(s_t)$

 Observe reward r_t and next state s_{t+1}

 Compute $\text{UCB}_{t+1}(g^{\pi})$ for all π based on $\text{UCB}_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $\text{UCB}_t(g^{\pi})$ with unknown dynamics? ✓

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

Gain Optimism: Implementation

Tentative algorithm

Observe state s_1

for $t = 1, 2, \dots$ **do**

Compute $\pi_t \leftarrow \arg \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\}$


Take action $a_t = \pi_t(s_t)$

Observe reward r_t and next state s_{t+1}

Compute $\text{UCB}_{t+1}(g^{\pi})$ for all π based on $\text{UCB}_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

 *3 major issues:*

■ *Upper confidence bounds*: construct $\text{UCB}_t(g^{\pi})$ with unknown dynamics? 

■ How to efficiently *compute* $\max_{M \in \mathcal{M}_t} g_M^{\pi}$ for every π ?

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

Gain Optimism: Implementation

Tentative algorithm

Observe state s_1

for $t = 1, 2, \dots$ **do**

Compute $\pi_t \leftarrow \arg \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\}$


Take action $a_t = \pi_t(s_t)$

Observe reward r_t and next state s_{t+1}

Compute $\text{UCB}_{t+1}(g^{\pi})$ for all π based on $\text{UCB}_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

 *3 major issues:*

■ *Upper confidence bounds*: construct $\text{UCB}_t(g^{\pi})$ with unknown dynamics? 

■ How to efficiently *compute* $\max_{M \in \mathcal{M}_t} g_M^{\pi}$ for every π ?

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

Extended MDP

[Strehl and Littman, 2008, Jaksch et al., 2010]

Theorem (Bounded parameter MDP \iff Extended MDP)

Let $\mathcal{M}_t^+ := \langle \mathcal{S}, \mathcal{A}_t^+, r^+, p^+ \rangle$ be an *extended* MDP such that

$$\mathcal{A}_t^+(s) = \mathcal{A}(s) \times B_t^r(s, a) \times B_t^p(s, a)$$

with $a^+ = (a, r, p) \in \mathcal{A}_t^+(s)$, $r^+(s, a^+) = r$, $p^+(\cdot | s, a^+) = p$.

Continuous **compact**
action space

Then the optimal gain of \mathcal{M}_t^+ satisfies

$$g_{\mathcal{M}_t^+}^* := \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\}$$

Let $\pi_t^+ = \arg \max_{\pi} g_{\mathcal{M}_t^+}^{\pi}$, then

$$\pi_t = \arg \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\} \text{ s.t. } \pi_t(s) = \pi_t^+(s)[a]$$

Extended MDP

[Strehl and Littman, 2008, Jaksch et al., 2010]

Theorem (Bounded parameter MDP \iff Extended MDP)

Let $\mathcal{M}_t^+ := \langle \mathcal{S}, \mathcal{A}_t^+, r^+, p^+ \rangle$ be an *extended* MDP such that

$$\mathcal{A}_t^+(s) = \mathcal{A}(s) \times B_t^r(s, a) \times B_t^p(s, a)$$

with $a^+ =$ *Abuse of notation:* \mathcal{M}_t denotes the extended MDP compact space

Then the optimal gain of \mathcal{M}_t^+ satisfies

$$g_{\mathcal{M}_t^+}^* := \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\}$$

Let $\pi_t^+ = \arg \max_{\pi} g_{\mathcal{M}_t^+}^{\pi}$, then

$$\pi_t = \arg \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\} \text{ s.t. } \pi_t(s) = \pi_t^+(s)[a]$$

Extended Value Iteration

Value iteration on \mathcal{M}_t

$$\begin{aligned}
 v_{n+1}(s) &= \mathcal{L}_t v_n(s) = \max_{(a,r,p) \in \mathcal{A}(s) \times B_t^r(s,a) \times B_t^p(s,a)} \left\{ r + p^\top v_n \right\} \\
 &= \max_{a \in \mathcal{A}(s)} \left\{ \max_{r \in B_t^r(s,a)} r + \max_{p \in B_t^p(s,a)} p^\top v_n \right\} \\
 &= \max_{a \in \mathcal{A}(s)} \left\{ \hat{r}_t(s,a) + \beta_t^r(s,a) + \max_{p \in B_t^p(s,a)} p^\top v_n \right\}
 \end{aligned}$$

$\pi_t = \text{Greedy policy w.r.t. } v_n$

Gain Optimism: Implementation

Tentative algorithm

Observe state s_1

for $t = 1, 2, \dots$ do

Compute $\pi_t \leftarrow \arg \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\}$

Take action $a_t = \pi_t(s_t)$

Observe reward r_t and next state s_{t+1}

Compute $\text{UCB}_{t+1}(g^{\pi})$ for all π based on $\text{UCB}_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

 *3 major issues:*

■ *Upper confidence bounds*: construct $\text{UCB}_t(g^{\pi})$ with unknown dynamics ✓

■ How to efficiently *compute* $\max_{M \in \mathcal{M}_t} g_M^{\pi}$ for every π ? ✓

■ *Computational complexity*: exponential number of policies ✓

■ *Frequent policy update*: inefficient exploration

Gain Optimism: Implementation

Tentative algorithm

Observe state s_1

for $t = 1, 2, \dots$ **do**

Compute $\pi_t \leftarrow \arg \max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^{\pi} \right\}$

Take action $a_t = \pi_t(s_t)$

Observe reward r_t and next state s_{t+1}

Compute $\text{UCB}_{t+1}(g^{\pi})$ for all π based on $\text{UCB}_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

⚠ 3 major issues:

- *Upper confidence bounds*: construct $\text{UCB}_t(g^{\pi})$ with unknown dynamics ✓
 - How to efficiently *compute* $\max_{M \in \mathcal{M}_t} g_M^{\pi}$ for every π ? ✓
- *Computational complexity*: exponential number of policies ✓
- *Frequent policy update*: inefficient exploration

Optimism: Frequency of Policy Updates

Proposition [Ortner, 2010]

There exists an MDP s.t.

$\Omega(T)$ number of policy updates \implies *linear regret*.

\implies $o(T)$ number of policy updates

Final Algorithm: UCRL2

Initialize $t \leftarrow 1$

Observe state s_1

Initialize empirical means $\hat{r}_1 = r_{\max}$ and $\hat{p}_1 = (1/S, \dots, 1/S)^\top$

Initialize visit counts $N_1 = 0$

for *episodes* $k = 1, 2, \dots$ **do**

 Set $t_k \leftarrow t$

 Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$

 Using EVI, compute *optimistic policy* π_k and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\mathcal{M}_k}^{\pi_k} h_k = h_k + g_k e \quad \text{with} \quad g_k = g_{\mathcal{M}_k}^* \geq g_{M^*}^*$$

while $N_t(s_t, a_t) < \max\{1, N_{t_k}(s_t, a_t)\}$ **do**

 Take action $a_t = \pi_k(s_t)$

 Observe reward r_t and next state s_{t+1}

 Compute new empirical means $\hat{r}_{t+1}(s_t, a_t)$ and $\hat{p}_{t+1}(\cdot | s_t, a_t)$

 Compute new visit count $N_{t+1}(s_t, a_t)$

$t \leftarrow t + 1$

end

end

Final Algorithm: UCRL2

Initialize $t \leftarrow 1$

Observe state s_1

Initialize empirical means $\hat{r}_1 = r_{\max}$ and $\hat{p}_1 = (1/S, \dots, 1/S)^\top$

Initialize visit counts $N_1 = 0$

for *episodes* $k = 1, 2, \dots$ do

Set $t_k \leftarrow t$

Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$

Using EVI, compute *optimistic policy* π_k and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\mathcal{M}_k}^{\pi_k} h_k = h_k + g_k e \quad \text{with} \quad g_k = g_{\mathcal{M}_k}^* \geq g_{M^*}^*$$

while $N_t(s_t, a_t) < \max\{1, N_{t_k}(s_t, a_t)\}$ do

Take action $a_t = \pi_k(s_t)$

Observe reward r_t and next state s_{t+1}

Compute new empirical means $\hat{r}_{t+1}(s_t, a_t)$ and $\hat{p}_{t+1}(\cdot | s_t, a_t)$

Compute new visit count $N_{t+1}(s_t, a_t)$

$t \leftarrow t + 1$

end

end

Optimism

Final Algorithm: UCRL2

Initialize $t \leftarrow 1$

Observe state s_1

Initialize empirical means $\hat{r}_1 = r_{\max}$ and $\hat{p}_1 = (1/S, \dots, 1/S)^\top$

Initialize visit counts $N_1 = 0$

for *episodes* $k = 1, 2, \dots$ do

Set $t_k \leftarrow t$

Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$

Using EVI, compute *optimistic policy* π_k and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\pi_k}^{\mathcal{M}_k} h_k = h_k + g_k e \quad \text{with} \quad g_k = g_{\mathcal{M}_k}^* \geq g_{M^*}^*$$

Bellman equation in \mathcal{M}_k

while $N_t(s_t, a_t) < \max\{1, N_{t_k}(s_t, a_t)\}$ do

Take action $a_t = \pi_k(s_t)$

Observe reward r_t and next state s_{t+1}

Compute new empirical means $\hat{r}_{t+1}(s_t, a_t)$ and $\hat{p}_{t+1}(\cdot | s_t, a_t)$

Compute new visit count $N_{t+1}(s_t, a_t)$

$t \leftarrow t + 1$

end

end

Optimism

Final Algorithm: UCRL2

Initialize $t \leftarrow 1$

Observe state s_1

Initialize empirical means $\hat{r}_1 = r_{\max}$ and $\hat{p}_1 = (1/S, \dots, 1/S)^\top$

Initialize visit counts $N_1 = 0$

for *episodes* $k = 1, 2, \dots$ do

Set $t_k \leftarrow t$

Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$

Using EVI, compute *optimistic policy* π_k and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\pi_k}^{\mathcal{M}_k} h_k = h_k + g_k e \quad \text{with} \quad g_k = g_{\mathcal{M}_k}^* \geq g_{M^*}^*$$

Bellman equation in \mathcal{M}_k

while $N_t(s_t, a_t) < \max\{1, N_{t_k}(s_t, a_t)\}$ do

Take action $a_t = \pi_k(s_t)$

Observe reward r_t and next state s_{t+1}

Compute new empirical means $\hat{r}_{t+1}(s_t, a_t)$ and $\hat{p}_{t+1}(\cdot | s_t, a_t)$

Compute new visit count $N_{t+1}(s_t, a_t)$

$t \leftarrow t + 1$

end

end

Optimism

Stopping condition of an episode

UCRL2: Regret Guarantees

Theorem (Thm.2 of [Jaksch et al., 2010])

There exists a numerical constant $\beta > 0$ such that in any *communicating* MDP $M^* = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, with probability *at least* $1 - \delta$, UCRL2 suffers a regret bounded as

$$\forall T \geq 1, R(T, M^*, \text{UCRL2}) \leq \beta \cdot r_{\max} D S \sqrt{AT \log \left(\frac{T}{\delta} \right)}$$

UCRL2: Regret Guarantees

Theorem (Thm.2 of [Jaksch et al., 2010])

There exists a numerical constant $\beta > 0$ such that in any *communicating* MDP $M^* = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, with probability *at least* $1 - \delta$, UCRL2 suffers a regret bounded as

$$\forall T \geq 1, R(T, M^*, \text{UCRL2}) \leq \beta \cdot r_{\max} DS \sqrt{AT \log \left(\frac{T}{\delta} \right)}$$

Comparison to lower bound

$$\bar{R}(T, M^*, \text{UCRL}) \geq 0.015 \sqrt{DSAT}$$

UCRL2: Regret Guarantees

Theorem (Thm.2 of [Jaksch et al., 2010])

There exists a numerical constant $\beta > 0$ such that in any *communicating* MDP $M^* = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, with probability *at least* $1 - \delta$, UCRL2 suffers a regret bounded as

$$\forall T \geq 1, R(T, M^*, \text{UCRL2}) \leq \beta \cdot r_{\max} \textcolor{red}{DS} \sqrt{AT \log \left(\frac{T}{\delta} \right)}$$

Comparison to lower bound

$$\overline{R}(T, M^*, \text{UCRL}) \geq 0.015 \sqrt{\textcolor{red}{DS} AT}$$

- Can the gap between upper and lower bound be closed? [👉 More on this later](#)

UCRL2: Regret Guarantees (cont'd.)

Theorem (Thm.4 of [Jaksch et al., 2010])

There exists a numerical constant $\beta > 0$ such that in any *ergodic* MDP $M^* = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, for all $T \geq 1$, UCRL2 (with $\delta = 1/T$) suffers a regret bounded as

$$\overline{R}(T, M^*, \text{UCRL2}) \leq \beta \cdot r_{\max} \frac{D^2 S^2 A \log(T)}{\delta_g^*} + \text{Big constant independent of } T$$

with

$$\delta_g^* := g_{M^*}^* - \max_{s \in \mathcal{S}, \pi} \left\{ g_{M^*}^\pi(s) < g_M^* \right\} \sim \text{"gap in gain"}$$

UCRL2: Regret Guarantees (cont'd.)

Theorem (Thm.4 of [Jaksch et al., 2010])

There exists a numerical constant $\beta > 0$ such that in any *ergodic* MDP $M^* = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, for all $T \geq 1$, UCRL2 (with $\delta = 1/T$) suffers a regret bounded as

$$\overline{R}(T, M^*, \text{UCRL2}) \leq \beta \cdot r_{\max} \frac{D^2 S^2 A \log(T)}{\delta_g^*} + \text{Big constant independent of } T$$

with

$$\delta_g^* := g_{M^*}^* - \max_{s \in \mathcal{S}, \pi} \left\{ g_{M^*}^\pi(s) < g_{M^*}^* \right\} \sim \text{"gap in gain"}$$

Comparison to lower bound

$$\liminf_{T \rightarrow \infty} \frac{\overline{R}(T, M^*, \mathfrak{A})}{\log T} \geq K_{M^*}, \text{ with } K_{M^*} \lesssim \frac{D^2 S A}{\min_{s,a} \delta_{M^*}^*(s, a)}$$

UCRL2: Regret Guarantees (cont'd.)

Theorem (Thm.4 of [Jaksch et al., 2010])

There exists a numerical constant $\beta > 0$ such that in any *ergodic* MDP $M^* = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, for all $T \geq 1$, UCRL2 (with $\delta = 1/T$) suffers a regret bounded as

$$\bar{R}(T, M^*, \text{UCRL2}) \leq \beta \cdot r_{\max} \frac{D^2 \textcolor{red}{S}^2 A \log(T)}{\delta_g^*} + \text{Big constant independent of } T$$

with

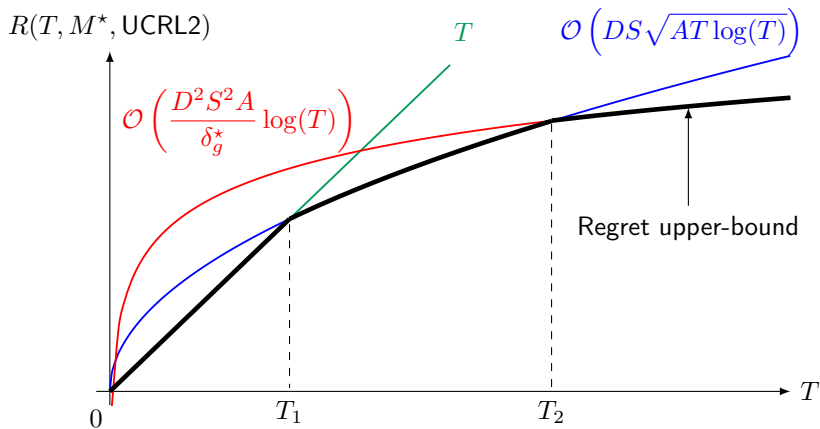
$\blacksquare \delta_g^* := g_{M^*}^* - \max_{s \in \mathcal{S}, \pi} \{g_{M^*}^\pi(s) < g_M^*\} \sim \text{“gap in gain”}$

how do they compare?

Comparison to lower bound

$$\liminf_{T \rightarrow \infty} \frac{\bar{R}(T, M^*, \mathfrak{A})}{\log T} \geq K_{M^*}, \text{ with } K_{M^*} \lesssim \frac{D^2 \textcolor{red}{S} A}{\min_{s,a} \delta_{M^*}^*(s,a)}$$

Qualitative Regret Shape



*illustrative plot

Regret Bound of UCRL2: Proof Sketch

$$1 \quad R(T, M^*, \text{UCRL2}) = \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} g_{M^*}^* - r(s_t, a_t) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} g_k - r(s_t, a_t)$$

Split in episodes

Optimism: $g_k \geq g_{M^*}^*$

Regret Bound of UCRL2: Proof Sketch

$$1 \quad R(T, M^*, \text{UCRL2}) = \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} g_{M^*}^* - r(s_t, a_t) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} g_k - r(s_t, a_t)$$

$$2 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} g_k = \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} r_k(s_t, a_t) + p_k(\cdot | s_t, a_t)^\top h_k - h_k(s_t)$$


Bellman equation ($a_t = \pi_k(s_t)$):

$$L_{\mathcal{M}_k}^{\pi_k} h_k(s_t) = h_k(s_t) + g_k$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} r_k(s_t, a_t) - r(s_t, a_t) + p_k(\cdot | s_t, a_t)^\top h_k - h_k(s_t)$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} r_k(s_t, a_t) - r(s_t, a_t) + p_k(\cdot | s_t, a_t)^\top h_k - h_k(s_t)$$


Assumption: true reward is known $r = r_k$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot | s_t, a_t)^\top h_k - h_k(s_t)$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot | s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{3} \quad p_k(\cdot | s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot | s_t, a_t) - p(\cdot | s_t, a_t) \right)^\top h_k + p(\cdot | s_t, a_t)^\top h_k - h_k(s_t)$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{3} \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\begin{aligned} \text{4} \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) &= \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_{t+1}) \\ &\quad + \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t) \end{aligned}$$

Regret Bound of UCRL2: Proof Sketch

$$1-2 \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$3 \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$4 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \Rightarrow \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_{t+1})$$

Martingale Difference Sequence
(Azuma's inequality)

$$+ \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t)$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{3} \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\begin{aligned} \text{4} \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) &\lesssim \sup_k \{ \text{sp}(h_k) \} \sqrt{T \log(T/\delta)} \\ &\quad + \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t) \end{aligned}$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{3} \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{4} \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim \sup_k \{ \text{sp}(h_k) \} \sqrt{T \log(T/\delta)}$$

$$+ \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t) \leftarrow \text{Telescopic sum}$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{3} \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{4} \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim \sup_k \{ \text{sp}(h_k) \} \sqrt{T \log(T/\delta)}$$

$$+ \textcolor{red}{m} \sup_k \{ \text{sp}(h_k) \}$$

Number of episodes
(stopping condition)

Regret Bound of UCRL2: Proof Sketch

$$1-2 \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$3 \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$4 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim \sup_k \{ \text{sp}(h_k) \} \sqrt{T \log(T/\delta)} \\ + SA \log(T) \sup_k \{ \text{sp}(h_k) \}$$

Regret Bound of UCRL2: Proof Sketch

$$1-2 \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$3 \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$4 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim \sup_k \{ \text{sp}(h_k) \} \sqrt{T \log(T/\delta)} \\ + SA \log(T) \sup_k \{ \text{sp}(h_k) \}$$

$$\text{sp}(h_k) \leq r_{\max} D \quad [\text{Bartlett and Tewari, 2009, Jaksch et al., 2010}]$$

Regret Bound of UCRL2: Proof Sketch

$$1-2 \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$3 \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$4 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log(T/\delta)} + r_{\max} D S A \log(T)$$

$$\text{sp}(h_k) \leq r_{\max} D \quad [\text{Bartlett and Tewari, 2009, Jaksch et al., 2010}]$$

Regret Bound of UCRL2: Proof Sketch

$$1-2 \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$3 \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$4 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log(T/\delta)} + r_{\max} D S A \log(T)$$

$$5 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k = \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} \underbrace{\left(p_k(\cdot|s_t, a_t) - \hat{p}_k(\cdot|s_t, a_t) \right)^\top h_k}_{\leq \text{sp}(h_k) \beta_k^p(s, a)} + \underbrace{\left(\hat{p}_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k}_{\leq \text{sp}(h_k) \beta_k^p(s, a)}$$

Regret Bound of UCRL2: Proof Sketch

$$\text{1-2} \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{3} \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$\text{4} \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log(T/\delta)} + r_{\max} D S A \log(T)$$

$$\text{5} \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k \lesssim r_{\max} D \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} \sqrt{\frac{S \log(T/\delta)}{N_{t_k}(s_t, a_t)}}$$

Regret Bound of UCRL2: Proof Sketch

$$1-2 \quad R(T, M^*, \text{UCRL2}) \leq \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$3 \quad p_k(\cdot|s_t, a_t)^\top h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k + p(\cdot|s_t, a_t)^\top h_k - h_k(s_t)$$

$$4 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\top h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log(T/\delta)} + r_{\max} D S A \log(T)$$

$$5 \quad \sum_{k=1}^m \sum_{t=t_k}^{t_{k+1}-1} \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\top h_k \lesssim r_{\max} D S \sqrt{A T \log(T/\delta)}$$

Refined Confidence Bounds

- UCRL2 with *Bernstein bounds* (instead of Hoeffding/Weissman):

 see [tutorial website](#)

$$R(T, M^*, \text{UCRL2B}) = \mathcal{O} \left(\sqrt{D \Gamma S A T \log \left(\frac{T}{\delta} \right) \log(T)} \right)$$

 Still not matching the lower bound!

 For most MPDs: $\Gamma \ll S$

Refined Confidence Bounds

- UCRL2 with *Bernstein bounds* (instead of Hoeffding/Weissman):

📖 see [tutorial website](#)

$$R(T, M^*, \text{UCRL2B}) = \mathcal{O} \left(\sqrt{D \Gamma S A T \log \left(\frac{T}{\delta} \right) \log(T)} \right)$$

🗨 Still not matching the lower bound!

👍 For most MPDs: $\Gamma \ll S$

- *Kullback-Leibler* UCRL [Filippi et al., 2010, Talebi and Maillard, 2018]:

$$R(T, M^*, \text{UCRL-KL}) = \mathcal{O} \left(\underbrace{\sqrt{\sum_{s,a} \mathbb{V}_{X \sim p^*(\cdot|s,a)} (h_{M^*}^*(X))}}_{\leq D^2 S A} S T \log \left(\frac{T}{\delta} \right) + D \sqrt{T} \right)$$

🗨 Only for ergodic MDPs!

Infinite Diameter (weakly communicating MDPs)

- *Known* bound on the optimal bias span $C \geq \text{sp}(h_{M^*}^*)$

[Bartlett and Tewari, 2009, Fruit et al., 2018b]

$$R(T, M^*, \text{SCAL}) = \mathcal{O} \left(\sqrt{C \Gamma S A T \log \left(\frac{T}{\delta} \right) \log(T)} \right)$$

🗨 Requires prior knowledge!

Infinite Diameter (weakly communicating MDPs)

- *Known* bound on the optimal bias span $C \geq \text{sp}(h_{M^*}^*)$

[Bartlett and Tewari, 2009, Fruit et al., 2018b]

$$R(T, M^*, \text{SCAL}) = \mathcal{O} \left(\sqrt{C \Gamma S A T \log \left(\frac{T}{\delta} \right) \log(T)} \right)$$

🗨 Requires prior knowledge!

- No prior knowledge: TUCRL [Fruit et al., 2018a]:

$$R(T, M^*, \text{SCAL}) = \mathcal{O} \left(\sqrt{D_{\text{com}} S_{\text{com}} \Gamma A T \log \left(\frac{T}{\delta} \right) \log(T)} \right)$$

🗨 Never achieves *logarithmic* regret! Intrinsic limitation of the setting!

Open Questions

1 *Tightness of minimax $\mathcal{O}(\sqrt{T})$ regret bounds for infinite horizon problems*

- Dependency on T : regret + sample complexity bounds?
- Analysis not tight *vs.* change in the algorithm?
- Lower bound not tight?

2 *Finite time logarithmic upper and lower regret bounds*

- Non-asymptotic lower bounds
- Tighter analysis of UCRL-like algorithms? New algorithms?

- 1 Setting the Stage
- 2 Lower Bounds
- 3 Optimism in Face of Uncertainty
- 4 Posterior Sampling
- 5 Asymptotically Optimal Algorithms
- 6 Extensions and Other Settings
- 7 Conclusion

Posterior Sampling

a.k.a. Thompson Sampling [Thompson, 1933]

Keep Bayesian posterior for the *unknown* MDP

👍 A sample from the posterior is used as an estimate of the unknown MDP

Exploration

Few samples \Rightarrow uncertainty in the estimate

More samples \Rightarrow posterior concentrates on the true MDP

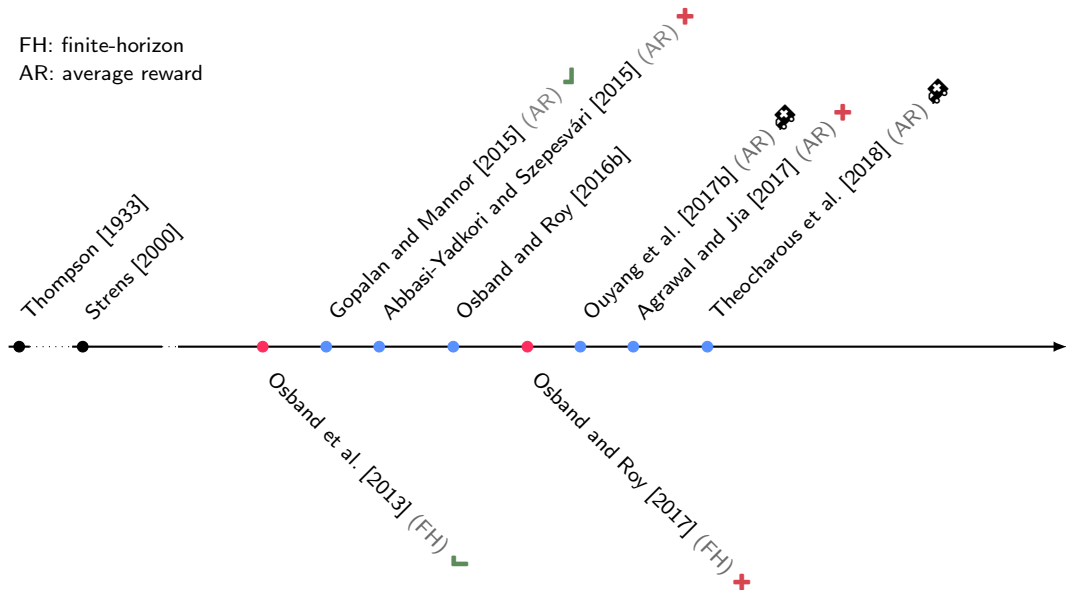
Exploitation

Set of MDPs



History: PS for Regret Minimization in RL

FH: finite-horizon
AR: average reward



Posterior Sampling

```
t ← 1
for episode k = 1, 2, ... do
  t_k ← t
   $M_k \sim \mu_{t_k}$ 
   $\pi_k \in \arg \max_{\pi} \{g_{M_k}^{\pi}\}$ 
  while not enough knowledge do
    Take action  $a_t \sim \pi_k(\cdot | s_t)$ 
    Observe reward  $r_t$  and next state  $s_{t+1}$ 
    Compute  $\mu_{t+1}$  based on  $\mu_t$  and
       $(s_t, a_t, r_t, s_{t+1})$ 
    t ← t + 1
  end
end
```

Posterior Sampling

```

t ← 1
for episode k = 1, 2, ... do
  t_k ← t
   $M_k \sim \mu_{t_k}$ 
   $\pi_k \in \arg \max_{\pi} \{g_{M_k}^{\pi}\}$ 
  while not enough knowledge do
    Take action  $a_t \sim \pi_k(\cdot | s_t)$ 
    Observe reward  $r_t$  and next state  $s_{t+1}$ 
    Compute  $\mu_{t+1}$  based on  $\mu_t$  and  $(s_t, a_t, r_t, s_{t+1})$ 
    t ← t + 1
  end
end

```

Prior distribution:

$$\forall \Theta, \mathbb{P}(M^* \in \Theta) = \mu_1(\Theta)$$

Posterior distribution:

$$\forall \Theta, \mathbb{P}(M^* \in \Theta | H_t, \mu_1) = \mu_t(\Theta)$$

Priors

- Dirichlet (transitions)
- Beta, Normal-Gamma, etc. (rewards)

Bayesian Regret

$$R^B(T, \mu_1, \mathfrak{A}) = \mathbb{E}_{M^* \sim \mu_1} \left[\underbrace{\bar{R}(T, M^*, \mathfrak{A})}_{:= \mathbb{E}[R(T, M^*, \mathfrak{A})]} \right] = \mathbb{E} \left[\sum_{t=1}^T g_{M^*}^* - r(s_t, a_t) \right]$$

TSDE: Thompson Sampling with Dynamic Episodes

[Ouyang et al., 2017b]

Episode length $l_k = t_{k+1} - t_k$ is dynamically determined by

- 1 Doubling of visits (stochastic)
- 2 Increasing length of previous episode by one (deterministic)

$$t_{k+1} = \min \left\{ t > t_k : \underbrace{\exists(s, a), N_t(s, a) > 2N_{t_k}(s, a)}_{(ST1)} \text{ or } \underbrace{t > t_k + l_{k-1}}_{(ST2)} \right\}$$

👉 (ST2) is $\sigma(H_{t_k})$ -measurable

$$l_k \leq l_{k-1} + 1$$

TSDE: Regret Guarantees

Theorem ([Ouyang et al., 2017b])

There exists a numerical constant $\beta > 0$ such that for any prior μ_1 whose support is a subset of *communicating* MDPs, TSDE suffers a regret bounded as

$$\forall T \geq 1, \quad R^B(T, \mu_1, \text{TSDE}) \leq \beta \cdot \left(CS \sqrt{AT \log(AT)} \right)$$

where

$$\mu_1 \quad \text{is such that} \quad \sup_{M^* \sim \mu_1} \left\{ sp(h_{M^*}^*) \right\} \leq C < +\infty \quad (\text{ASM-SP})$$

Proof Step 1: Regret Decomposition

- 👉 The support of the prior μ_1 is a subset of communicating MDPs
 M_k is communicating and optimality equation (i.e., constant gain)

$$\begin{aligned}
 R^B(T, \mu_1, \text{TSDE}) &\leq \underbrace{T \mathbb{E}[g_{M^*}^*] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k g_{M_k}^* \right]}_{R_g} \\
 &+ \mathbb{E} \left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_t) - h_k(s_{t+1})) \right] \\
 &+ \mathbb{E} \left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_{t+1}) - p_k(\cdot | s_t, a_t)^\top h_k) + r_k(s_t, a_t) - r(s_t, a_t) \right]
 \end{aligned}$$

Proof Step 1: Regret Decomposition

- 👉 The support of the prior μ_1 is a subset of communicating MDPs
 M_k is communicating and optimality equation (i.e., constant gain)

$$\begin{aligned}
 R^B(T, \mu_1, \text{TSDE}) &\leq \underbrace{T \mathbb{E}[g_{M^*}^*] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k g_{M_k}^* \right]}_{R_g} \\
 &+ \mathbb{E} \left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_t) - h_k(s_{t+1})) \right] \\
 &+ \mathbb{E} \left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_{t+1}) - p_k(\cdot | s_t, a_t)^\top h_k) + r_k(s_t, a_t) - r(s_t, a_t) \right]
 \end{aligned}$$

Telescopic sum
+ span bound (ASM-SP)[†]

[†] as in UCRL2

Proof Step 1: Regret Decomposition

- 👉 The support of the prior μ_1 is a subset of communicating MDPs
 M_k is communicating and optimality equation (i.e., constant gain)

$$\begin{aligned}
 R^B(T, \mu_1, \text{TSDE}) &\leq \underbrace{T \mathbb{E}[g_{M^*}^*] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k g_{M_k}^* \right]}_{R_g} \\
 &\quad + \mathbb{E} \left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_t) - h_k(s_{t+1})) \right] \\
 &\quad + \mathbb{E} \left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_{t+1}) - p_k(\cdot | s_t, a_t)^\top h_k) + r_k(s_t, a_t) - r(s_t, a_t) \right]
 \end{aligned}$$

Confidence sets[†] →

Telescopic sum
+ span bound (ASM-SP)[†] →

[†] as in UCRL2

Proof Step 2: Bounding R_g

Thompson Sampling Lemma [Osband et al., 2013, Ouyang et al., 2017b]

Let t_k be an almost surely finite $\sigma(H_{t_k})$ -stopping time. For any measurable function f and $\sigma(H_{t_k})$ -measurable variable X

$$\mathbb{E}[f(M_k, X)|H_{t_k}] = \mathbb{E}[f(M^*, X)|H_{t_k}]$$

Proof Step 2: Bounding R_g

$$R_g = \mathbb{E} \left[\sum_{t=1}^{k_T} l_t \ g_{M^*}^* \right] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k \ g_{M_k}^* \right]$$

Proof Step 2: Bounding R_g

random duration of episode k
not $\sigma(H_{t_k})$ -measurable

$$R_g = \mathbb{E} \left[\sum_{t=1}^{k_T} l_t \ g_{M^*}^* \right] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k \ g_{M_k}^* \right]$$

Proof Step 2: Bounding R_g

random duration of episode k
not $\sigma(H_{t_k})$ -measurable

$$\begin{aligned}
 R_g &= \mathbb{E} \left[\sum_{t=1}^{k_T} l_t g_{M^*}^* \right] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k g_{M_k}^* \right] \\
 &\leq \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1) \left(g_{M^*}^* - g_{M_k}^* \right) \right] + \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1 - l_k) g_{M_k}^* \right] \quad \left(\begin{array}{l} \text{by (ST2)} \\ l_k \leq l_{k-1} + 1 \end{array} \right)
 \end{aligned}$$

Proof Step 2: Bounding R_g

random duration of episode k
not $\sigma(H_{t_k})$ -measurable

$$\begin{aligned}
 R_g &= \mathbb{E} \left[\sum_{t=1}^{k_T} l_k g_{M^*}^* \right] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k g_{M_k}^* \right] \\
 &\leq \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1) \left(g_{M^*}^* - g_{M_k}^* \right) \right] + \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1 - l_k) g_{M_k}^* \right] \quad \left(\begin{array}{l} \text{by (ST2)} \\ l_k \leq l_{k-1} + 1 \end{array} \right) \\
 &\leq \quad \quad \quad 0 \quad \quad \quad + r_{\max} \mathbb{E}[k_T]
 \end{aligned}$$

t_k is a stopping time
 $(l_{k-1} + 1)$ is $\sigma(H_{t_k})$ -measurable
 \implies use TS lemma

Proof Step 2: Bounding R_g

random duration of episode k
not $\sigma(H_{t_k})$ -measurable

$$\begin{aligned}
 R_g &= \mathbb{E} \left[\sum_{t=1}^{k_T} l_k g_{M^*}^* \right] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k g_{M_k}^* \right] \\
 &\leq \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1) (g_{M^*}^* - g_{M_k}^*) \right] + \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1 - l_k) g_{M_k}^* \right] \quad \left(\begin{array}{l} \text{by (ST2)} \\ l_k \leq l_{k-1} + 1 \end{array} \right) \\
 &\leq \quad \quad \quad 0 \quad \quad \quad + r_{\max} \mathbb{E}[k_T]
 \end{aligned}$$

t_k is a stopping time
 $(l_{k-1} + 1)$ is $\sigma(H_{t_k})$ -measurable
 \implies use TS lemma

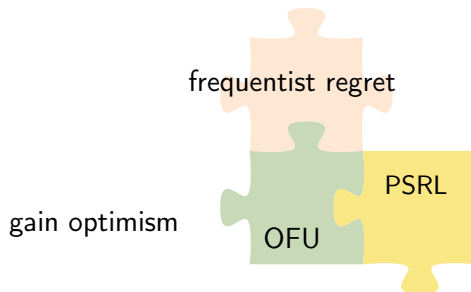
$$\begin{aligned}
 \sum_{k=1}^{k_T} l_{k-1} &= l_0 + \sum_{k=1}^{k_T-1} l_k \leq T \\
 g_{M_k}^* &\in [0, r_{\max}], \forall k
 \end{aligned}$$

$$\begin{aligned}
R_g &= \mathbb{E} \left[\sum_{t=1}^{k_T} l_t g_{M^*}^* \right] - \mathbb{E} \left[\sum_{k=1}^{k_T} l_k g_{M_k}^* \right] \\
&\leq \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1) (g_{M^*}^* - g_{M_k}^*) \right] + \mathbb{E} \left[\sum_{k=1}^{k_T} (l_{k-1} + 1 - l_k) g_{M_k}^* \right] \quad \left(\begin{array}{c} \text{by (ST2)} \\ l_k \leq l_{k-1} + 1 \end{array} \right) \\
&\leq 0 + r_{\max} \mathbb{E}[k_T]
\end{aligned}$$
$$\begin{array}{l} \leq \\ \leq r_{\max} \sqrt{2SAT \log(T)} \end{array} \quad \begin{array}{c} \uparrow \\ 0 \end{array} \quad + \quad \begin{array}{c} \uparrow \\ r_{\max} \mathbb{E}[k_T] \end{array} \quad ([\text{Ouyang et al., 2017b}] \text{ Lem. 1})$$
$$\sum_{k=1}^{k_T} l_{k-1} = l_0 + \sum_{k=1}^{k_T-1} l_k \leq T$$

$$g_{M_k}^{\pi_k} \in [0, r_{\max}], \forall k$$

OPT-PSRL: Optimistic Posterior Sampling

[Agrawal and Jia, 2017]



1. Sample posterior $\psi = \tilde{O}(S)$ times

$$p_{sa}^i \sim \mu_{t_k}(s, a), \quad i = 1, \dots, \psi$$



\mathcal{M}_k is an *discrete extended* MDP

$$\tilde{p}(\cdot, s, a^i) = p_{s,a}^i, \quad a^i \in \mathcal{A} \times \{1, \dots, \psi\}$$

2. Solve \mathcal{M}_k for π_k

$$g_{M_k}^* \geq g_{M^*}^* - \tilde{O}\left(D\sqrt{SA/T}\right)$$

OPT-PSRL: Regret Guarantees

Theorem ([Agrawal and Jia, 2017])

There exists a numerical constant $\alpha, \beta > 0$ such that in any *communicating* MDP M^* , with probability *at least* $1 - \delta$ and for any $T \geq \alpha DA \log^2(T/\delta)$, Opt-PSRL suffers a regret bounded as:


$$R(T, M^*, \text{Opt-PSRL}) \leq \beta r_{\max} \cdot \left(DS \sqrt{AT \log \left(\frac{T}{\delta} \right)} + \text{poly}(S, A) DT^{1/4} \log \left(\frac{T}{\delta} \right) \right)$$

Open Questions

1 *The nature of bounded bias span assumption (Asm. ASM-SP)*

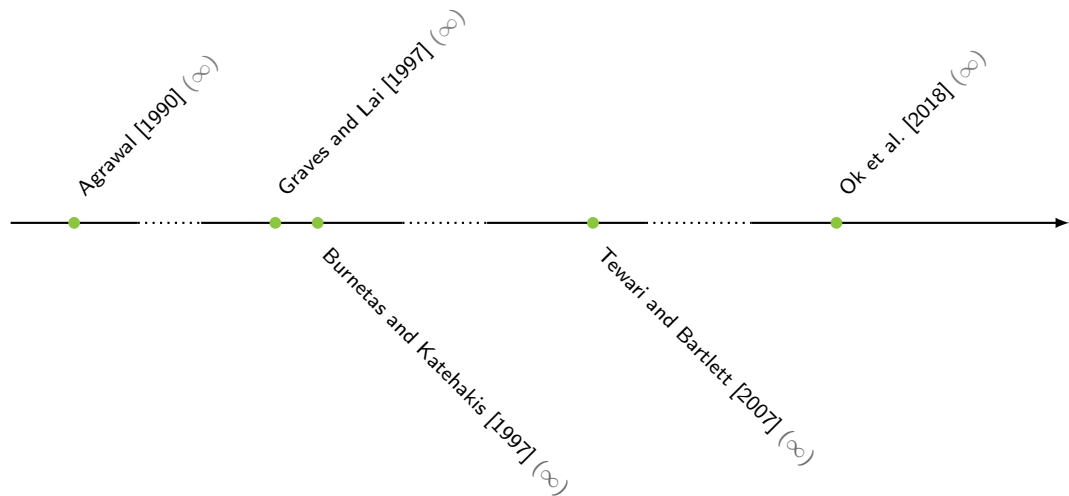
- Used in [Ouyang et al., 2017b, Theodorou et al., 2018]
- $\text{supp}(\mu_1)$ is continuous, then $\sup_{M^* \sim \mu_1} \{\text{sp}(h_{M^*}^*)\} = +\infty$ [e.g., Fruit et al. [2018a]]

2 *Statistical efficiency of PSRL*

- Claimed efficient Bayesian or frequentist $\tilde{O}(D\sqrt{SAT})$ regret bound
- Not supported by proofs, incorrect Lem. C.1 [Osband and Roy, 2016a] and Lem. C.2 [Agrawal and Jia, 2017] [ see tutorial website]

- 1 Setting the Stage
- 2 Lower Bounds
- 3 Optimism in Face of Uncertainty
- 4 Posterior Sampling
- 5 Asymptotically Optimal Algorithms
- 6 Extensions and Other Settings
- 7 Conclusion

History: Asymptotic Regret Minimization



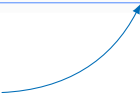
Asymptotic Lower-Bound

Theorem (Thm. 2, [Burnetas and Katehakis, 1997])

Any algorithm \mathfrak{A} s.t. $\overline{R}(T, M, \mathfrak{A}) = o(T^\alpha)$ for all $\alpha > 0$ and *ergodic* MDP M should satisfy

$$\forall (s, a) : \mathcal{M}_{M^*}^{alt}(s, a), \quad \liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_T(s, a)]}{\log T} \geq \frac{1}{\inf_{M \in \mathcal{M}_{M^*}^{alt}(s, a)} KL_{M^*, M}(s, a)}$$

👍 Should be satisfied by optimal algorithms
necessary to be uniformly good on all the possible *alternative* models



BKIA: Burnetas-Katehakis Index Algorithm

[Burnetas and Katehakis, 1997]

for $t = 1, \dots, T$ **do**

$D_t(s) \leftarrow \{a \in \mathcal{A}(s) : N_t(s, a) \geq \log^2(N_t(s))\}$
 $(g_t, h_t) \leftarrow \text{solve } \widehat{M}_t = \langle \mathcal{S}, D_t, \widehat{p}_t, r \rangle$

A Solve empirical MDP \widehat{M}_t on a restricted action set

if $\exists a \in \Pi_{\widehat{M}_t}^*(s_t), N_t(s_t, a) \geq \log^2(N_t(s_t) + 1)$ **then**

$a_t \in \arg \max_{a \in \mathcal{A}(s_t)} \{b_t(s, a; h_t)\}$

B Select maximum index action

else

$a_t \in \arg \min_{a \in \Pi_{\widehat{M}_t}^*(s_t)} \{N_t(s, a)\}$

C Force exploration of “underestimated” actions

end

Observe reward r_t and next state s_{t+1}

end

BKIA: Interpretation

B Exploration & Exploitation

$$a_t \in \arg \max_{a \in \mathcal{A}} \{b_t(s_t, a)\} \longrightarrow \oplus \longrightarrow \boxed{\text{Optimistic greedy}}$$

$$b_t(s, a) = \sup_{q \in \Delta(\mathcal{S})} \left\{ L_q^a h_{\widehat{M}_t}^*(s) : N_t(s, a) \text{KL}(\widehat{p}_t(\cdot | s_t, a) \| q) \leq \log(t) \right\}$$

$$\text{related to } - \inf_{M \in \mathcal{M}_{\widehat{M}_t}^{\text{alt}}(s, a)} \left\{ \delta_{\widehat{M}_t}^*(s, a) : N_t(s, a) \text{KL}_{\widehat{M}_t, M}(s, a) \leq \log(t) \right\}$$

⚠ A not so explicit way of controlling the lower bound

BKIA: Interpretation

B Exploration & Exploitation

$$a_t \in \arg \max_{a \in \mathcal{A}} \{b_t(s_t, a)\} \longrightarrow \oplus \longrightarrow \text{Optimistic greedy}$$

$$b_t(s, a) = \sup_{q \in \Delta(\mathcal{S})} \left\{ L_q^a h_{\widehat{M}_t}^*(s) : N_t(s, a) \text{KL}(\widehat{p}_t(\cdot | s_t, a) \| q) \leq \log(t) \right\}$$

related to $-\inf_{M \in \mathcal{M}_{\widehat{M}_t}^{\text{alt}}(s, a)} \left\{ \delta_{\widehat{M}_t}^*(s, a) : N_t(s, a) \text{KL}_{\widehat{M}_t, M}(s, a) \leq \log(t) \right\}$

⚠ A not so explicit way of controlling the lower bound

📖 Computing b_t is similar to KL-UCB [Garivier and Cappé, 2011] for MAB.

BKIA: Interpretation

Forced Exploration

when $\forall a \in \Pi_{\widehat{M}_t}^*(s_t), N_t(s_t, a) < \log^2(N_t(s_t) + 1)$

- BKIA prevents that *all* optimal actions *will become* under-explored

$$\implies a_t \in \Pi_{\widehat{M}_t}^*(s_t)$$

 Asymptotic monotonic property

$$\mathbb{P}\left(g_{M^*(D_{t+1})}^* \geq g_{M^*(D_t)}^*\right) = 1 - o\left(\frac{1}{t}\right) \quad \text{as } t \rightarrow \infty$$

BKIA: Regret Guarantees

Theorem (Thm. 1, [Burnetas and Katehakis, 1997])

For any *ergodic* MDP M^* , the expected regret of BKIA is upper bounded as

$$\limsup_{T \rightarrow \infty} \frac{\overline{R}(T, M^*, BKIA)}{\log T} \leq K_{M^*}^*$$

BKIA: Regret Guarantees

Theorem (Thm. 1, [Burnetas and Katehakis, 1997])

For any *ergodic* MDP M^* , the expected regret of BKIA is upper bounded as

$$\limsup_{T \rightarrow \infty} \frac{\overline{R}(T, M^*, BKIA)}{\log T} \leq K_{M^*}^*$$

👍 OLP [Tewari and Bartlett, 2007] replaces the KL constraint with an L_1

BKIA: Regret Proof

By [Prop. 1, [Burnetas and Katehakis, 1997]]

$$\overline{R}(T, M^*, \mathfrak{A}) = \sum_s \sum_{a \notin \Pi_{M^*}^*(s)} \mathbb{E}[N_T(s, a)] \delta_{M^*}^*(s, a) + O(1), \quad \text{as } T \rightarrow +\infty$$

We define W_T^1 s.t.

$$\mathbb{E}[N_T(s, a)] \leq \mathbb{E}[W_T^1(s, a, \varepsilon)] + o(\log T)$$

Ergodicity of MDP (g and h continuity)
about $h_{\widehat{M}_t}^* \rightarrow h_{M^*}^*$

BKIA: Regret Proof

Event

$$E_t^1 = \left\{ \|h_{\widehat{M}_t}^* - h_{M^*}^*\|_\infty \leq \varepsilon \wedge \Pi_{\widehat{M}_t}^*(s) \subseteq \Pi_{M^*}^*(s), \forall s \right\}$$

$$\widehat{M}_t \approx M^*$$

$$E_t^2 = \{b_t(s, a) < L_{M^*}^* h_{M^*}^*(s) - 2\varepsilon\}$$

BKIA: Regret Proof

Event

$$E_t^1 = \left\{ \|h_{\widehat{M}_t}^* - h_{M^*}^*\|_\infty \leq \varepsilon \wedge \Pi_{\widehat{M}_t}^*(s) \subseteq \Pi_{M^*}^*(s), \forall s \right\}$$

$$\widehat{M}_t \approx M^*$$

$$E_t^2 = \{b_t(s, a) < L_{M^*}^* h_{M^*}^*(s) - 2\varepsilon\}$$

$$W_T^1(s, a, \varepsilon) = \sum_{t=1}^T \mathbb{1}(s_t, a_t = s, a) \times \mathbb{1}(E_t^1 \wedge E_t^2)$$

? One Step Optimism

$$\forall (s, a) : \mathcal{M}_{M^*}^{\text{alt}}(s, a) \neq \emptyset$$

$$\lim_{\varepsilon \rightarrow 0} \limsup_{T \rightarrow \infty} \frac{\mathbb{E}[W_T^1(s, a, \varepsilon)]}{\log T} \leq \frac{1}{\inf_{M \in \mathcal{M}_{M^*}^{\text{alt}}(s, a)} \text{KL}_{M^*, M}(s, a)}$$

DEL: Directed Exploration Learning

[Ok et al., 2018]

- DEL exploits the same idea of BKIA

Explore suboptimal actions no more than what prescribed by the lower bound

- Exploration rate of sub-optimal action is *directed by the lower bound*

$$\text{target } \eta_t(s, a) \approx \mathbb{E}[N_T(s, a)]$$



OSSB [Combes et al., 2017] asymptotic optimal algorithm for structured bandit

for $t = 1, \dots, T$ do

$D_t(s) \leftarrow \{a \in \mathcal{A}(s) : N_t(s, a) \geq \log^2(N_t(s))\}$
 $(g_t, h_t) \leftarrow \text{solve } \widehat{M}_t = \langle S, D_t, \widehat{p}_t, r \rangle$

if $\forall a \in \Pi_{\widehat{M}_t}^*(s_t), N_t(s_t, a) < \log^2(N_t(s_t) + 1)$ then

$a_t \in \arg \min_{a \in \Pi_{\widehat{M}_t}^*(s_t)} \{N_t(s, a)\}$

else if $C^{xpt}(H_t)$ then

B1 *exploit* ($a_t \in \Pi_{\widehat{M}_t}^*(s_t)$)

else

B2 *explore*

end

Observe reward r_t and next state s_{t+1}

end

A Solve empirical MDP \widehat{M}_t on a restricted action set

C Force exploration of “underestimated” actions

! BKIA automatically trade-off exploration and exploitation
 $B1 + B2 \approx B_{\text{BKIA}}$

DEL: Exploration

B2 Directly *optimize the lower bound* on the estimated MDP \widehat{M}_t

$$\eta_t = \arg \inf_{\eta \in \mathbb{R}^{S \times A}} \sum_{s,a} \eta(s,a) \delta_{\widehat{M}_t}^*(s,a)$$

$$\text{s.t. } \sum_{s,a} \eta(s,a) \text{KL}_{\widehat{M}_t, M}(s,a) \geq 1 \quad \forall M \in \mathcal{M}_{\widehat{M}_t}^{\text{alt}}(s,a)$$

$$a_t \in \arg \min_{\mathcal{A}: N_t(s_t, a) \leq \eta_t(s_t, a) \gamma_t} \{N_t(s_t, a)\} \quad * \gamma_t = (1 + \gamma)(1 + \log t)$$

DEL: Exploration

B2 Directly *optimize the lower bound* on the estimated MDP \widehat{M}_t

$$\begin{aligned} \eta_t &= \arg \inf_{\eta \in \mathbb{R}^{S \times A}} \sum_{s,a} \eta(s,a) \delta_{\widehat{M}_t}^*(s,a) \\ \text{s.t. } \sum_{s,a} \eta(s,a) \text{KL}_{\widehat{M}_t, M}(s,a) &\geq 1 \quad \forall M \in \mathcal{M}_{\widehat{M}_t}^{\text{alt}}(s,a) \end{aligned}$$

$$a_t \in \arg \min_{A: N_t(s_t, a) \leq \eta_t(s_t, a) \gamma_t} \{N_t(s_t, a)\} \quad * \gamma_t = (1 + \gamma)(1 + \log t)$$

💡 Lower bound sets the desired number of visits

$$\eta_t(s_t, a) \approx \mathbb{E}_{\widehat{M}_t} \left[N_T(s_t, a) \right] \approx \mathbb{E}_{M^*} \left[N_T(s_t, a) \right]$$

then track it (in one step)

🗨️ η_t computed on \widehat{M}_t and not M^* (wrong target)

DEL: Regret Guarantees

Theorem (Thm. 4, [Ok et al., 2018])

For any *ergodic* MDP M^\star and under some technical conditions, for any $\gamma > 0$, the expected regret of $\text{DEL}(\gamma)$ is upper bounded as

$$\limsup_{T \rightarrow \infty} \frac{\overline{R}(T, M^\star, \text{DEL}(\gamma))}{\log T} \leq (1 + \gamma) K_{M^\star}^\star$$

DEL: Regret Guarantees

Theorem (Thm. 4, [Ok et al., 2018])

For any *ergodic* MDP M^* and under some technical conditions, for any $\gamma > 0$, the expected regret of $\text{DEL}(\gamma)$ is upper bounded as

$$\limsup_{T \rightarrow \infty} \frac{\overline{R}(T, M^*, \text{DEL}(\gamma))}{\log T} \leq (1 + \gamma) K_{M^*}^*$$

👍 DEL works for MDPs with structure (e.g., Lipschitz continuity)

Open Questions

■ *The role of forced exploration*

- Why do we need to force exploration?
- Is it due to the lack of long-term optimism?
- Is it really required at algorithmic level?

■ *Finite Time Analysis*

■ *Refined lower bound*

- Current lower bound is derived from a bandit perspective

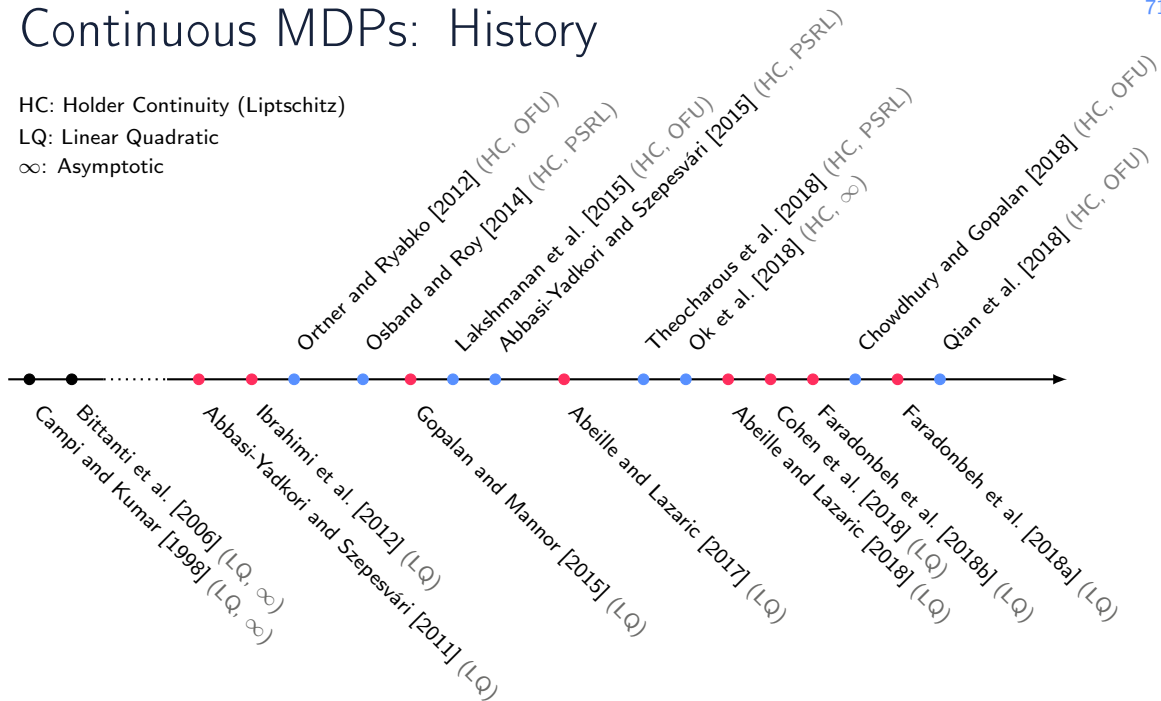
- 1 Setting the Stage
- 2 Lower Bounds
- 3 Optimism in Face of Uncertainty
- 4 Posterior Sampling
- 5 Asymptotically Optimal Algorithms
- 6 Extensions and Other Settings**
- 7 Conclusion

Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space \mathcal{S} , $|\mathcal{S}| = S < \infty$
 - Action space \mathcal{A} , $|\mathcal{A}| = A < \infty$
- } **finite**
- Transition distribution $p(\cdot | s, a) \in \Delta(\mathcal{S})$
 - Reward distribution with expectation $r(s, a) \in [0, r_{\max}]$
- 👉 The process generates history $H_t = (s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t)$, with $s_{t+1} \sim p(\cdot | s_t, a_t)$

∞ : Asymptotic



Hölder Continuity

\mathcal{S} continuous
 \mathcal{A} discrete

$L, \alpha > 0$ s.t. $\forall s, s' \in \mathcal{S}, a \in \mathcal{A}$:

$$|r(s, a) - r(s', a)| \leq r_{\max} L |s - s'|^\alpha$$

$$\|p(\cdot|s, a) - p(\cdot|s', a)\|_1 \leq L |s - s'|^\alpha$$

HC1 Asm.

$$\text{sp}(h_{M^\star}^\star) \leq C$$

HC2 Asm.

Hölder Continuity

\mathcal{S} continuous
 \mathcal{A} discrete

$L, \alpha > 0$ s.t. $\forall s, s' \in \mathcal{S}, a \in \mathcal{A}$:

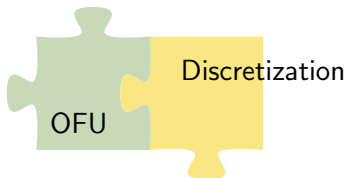
$$|r(s, a) - r(s', a)| \leq r_{\max} L |s - s'|^\alpha$$

$$\|p(\cdot | s, a) - p(\cdot | s', a)\|_1 \leq L |s - s'|^\alpha$$

HC1 Asm.

$$\text{sp}(h_{M^*}^*) \leq C$$

HC2 Asm.



[Ortner and Ryabko, 2012,
 Lakshmanan et al., 2015,
 Qian et al., 2018]

👍 L, α, C, T known in advance

OFU: Hölder Continuity

Theorem (Ortner and Ryabko [2012], Lakshmanan et al. [2015], Qian et al. [2018])

For any MDP M satisfying *Asm. (HC1) and (HC2)*, with probability at least $1 - \delta$ it holds that for any $T \geq 1$, the regret of UCCRL and SCCAL^+ is bounded as

$$R(T, M^*, \{\text{UCCRL}, \text{SCCAL}^+\}) \leq \beta \cdot CL \sqrt{A \log \left(\frac{T}{\delta} \right)} T^{(2+\alpha)/(2+2\alpha)}$$

If the transition function is *κ -times smoothly differentiable* ($\gamma = \alpha + \kappa$)

$$R(T, M^*, \text{UCCRL-KD}) \leq \beta \cdot CL \sqrt{A \log \left(\frac{T}{\delta} \right)} T^{(\gamma+2\alpha+\alpha\gamma)/(\gamma+\alpha+2\alpha\gamma)}$$

OFU: Lipschitz Continuity ($\alpha = 1$)

Theorem (Ok et al. [2018])

For any MDP M satisfying *Asm. (HC1) and (HC2)* with $\alpha = 1$ the regret of DEL is bounded as

$$\limsup_{T \rightarrow \infty} \frac{\bar{R}(T, M^*, \text{DEL})}{\log T} \leq S_L A \frac{(C+1)^3}{(\min_{s,a} \delta_{M^*}^*(s, a))^2}$$

with

$$S_L = \min\left\{S, \frac{8L(C+1)}{\min_{s,a} \delta_{M^*}^*(s, a)} + 1\right\}$$

Comparison

$$R(T, M^*, \{\text{UCCRL}, \text{SCCAL}^+\}) = \tilde{O}(T^{3/4})$$

$$R(T, M^*, \text{UCCRL-KD}) = \tilde{O}(T^{2/3}) \text{ as } \kappa \rightarrow \infty$$

Linear Quadratic Systems

$$\begin{aligned} \max_{\pi} \quad & \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right] \\ \text{s.t.} \quad & s_{t+1} = f(s_t, a_t, \epsilon_{t+1}) \\ & a_t \sim \pi(s_t) \end{aligned}$$

Linear Quadratic Systems

$$\max_{\pi} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T - \left(s_t^{\top} Q s_t + a_t^{\top} R a_t \right) \right]$$

$$\text{s.t.} \quad s_{t+1} = A s_t + B a_t + \epsilon_{t+1}$$

$$a_t \sim \pi(s_t)$$

Quadratic Reward

Linear Dynamics

LQ system $M = \langle A, B, Q, R \rangle$

Linear Quadratic Systems: Optimal Policy

■ *Optimal policy*

$$\pi_M^*(s) = \Sigma_M^* s$$

$$\Sigma_M^* = -(R + B^T P_M B)^{-1} (B^T P_M A)$$

solution of Discrete Algebraic
Riccati Equation (DARE)

■ *Optimal gain*

$$g_M^* = \text{Tr}(P_M)$$


Linear Quadratic Systems: Optimal Policy

■ *Optimal policy*

$$\pi_M^*(s) = \Sigma_M^* s$$

$$\Sigma_M^* = -(R + B^T P_M B)^{-1} (B^T P_M A)$$

solution of Discrete Algebraic
Riccati Equation (DARE)



■ *Optimal gain*

$$g_M^* = \text{Tr}(P_M)$$

if (A, B) are controllable, Σ_M^* makes the system *stable*

OFU-LQ

[Abbasi-Yadkori and Szepesvári, 2011]

assume Q and R are known


Optimism in LQ

■ Estimation

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle$$

Regularized Least Squares

where $(\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$



OFU-LQ

[Abbasi-Yadkori and Szepesvári, 2011]

assume Q and R are known

Optimism in LQ

■ Estimation

Statistically
admissible
models

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle$$

where

$$(\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$$

Regularized Least Squares

$$B_t^{\text{RLS}} = \{ \theta : \text{Tr}((\theta - \widehat{\theta}_t)^\top V_t (\theta - \widehat{\theta}_t)) \leq \beta_t \}$$

OFU-LQ

[Abbasi-Yadkori and Szepesvári, 2011]

assume Q and R are known

Optimism in LQ

■ Estimation

Statistically
admissible
models

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle$$

where $(\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$

Regularized Least Squares

$$B_t^{\text{RLS}} = \{ \theta : \text{Tr}((\theta - \widehat{\theta}_t)^\top V_t (\theta - \widehat{\theta}_t)) \leq \beta_t \}$$

■ Planning

$$\theta_t = \arg \max_{\theta \in \Theta \cap B_t^{\text{RLS}}} \{g_\theta^*\}$$

so that θ_t is
controllable

OFU-LQ

[Abbasi-Yadkori and Szepesvári, 2011]

assume Q and R are known

Optimism in LQ

■ Estimation

Statistically
admissible
models

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle$$

where

$$(\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$$

Regularized Least Squares

$$B_t^{\text{RLS}} = \{ \theta : \text{Tr}((\theta - \widehat{\theta}_t)^\top V_t (\theta - \widehat{\theta}_t)) \leq \beta_t \}$$

so that θ_t is
controllable

■ Planning

$$\theta_t = \arg \max_{\theta \in \Theta \cap B_t^{\text{RLS}}} \{g_\theta^*\}$$

🗨 Hard non-convex optimization problem

OFU-LQ: Regret

Theorem ([Abbasi-Yadkori and Szepesvári, 2011])

For any $\delta \in]0, 1[$, for any time T , with probability at least $1 - \delta$, the regret of OFU-LQ algorithm is bounded as

$$R(T, M^*, \text{OFU-LQ}) = \tilde{O}(\sqrt{T \log(1/\delta)})$$

OFU-LQ: Regret

Theorem ([Abbasi-Yadkori and Szepesvári, 2011])

For any $\delta \in]0, 1[$, for any time T , with probability at least $1 - \delta$, the regret of OFU-LQ algorithm is bounded as

$$R(T, M^*, \text{OFU-LQ}) = \tilde{O}(\sqrt{T \log(1/\delta)})$$

💡 major challenge

$\Sigma_{M^*}^* \rightarrow M^*$ stable controller ✓

$\Sigma_t \rightarrow M^*$???

central to the proof is how to control $\|s_t\|$

Open Question

Hölder continuity

1 *Posterior Sampling*

- [Theocharous et al., 2018] proved $\tilde{O}(C\sqrt{T})$
 - Under Asm. ASM-SP and Hölder continuity
 - Only for system parametrized by 1-dimensional parameter

2 Matching Lower Bound

LQ Systems

1 *Posterior Sampling*

- [Ouyang et al., 2017a] prove $\tilde{O}(\sqrt{T})$ Bayesian regret under restrictive assumptions
- [Abeille and Lazaric, 2017, 2018] proved $\tilde{O}(\sqrt{T})$ regret for PSRL with rejection sampling but only for 1-dimensional systems

2 *Efficient* OFU: many recent advances [Faradonbeh et al., 2018a, Cohen et al., 2019]

Other Settings

- Non-realizable approximated MDP (e.g. [Jiang et al., 2017])
- Non-stationary/adversarial environments (e.g. [Even-Dar et al., 2009, Neu et al., 2014])
- MDPs with arbitrary structure (e.g. [Gopalan and Mannor, 2015])
- Hierarchical exploration (e.g. [Fruit and Lazaric, 2017, Fruit et al., 2017])
- Low-exploration MDPs (e.g. [Zanette and Brunskill, 2018])
- Active/unsupervised exploration (e.g. [Lim and Auer, 2012, Hazan et al., 2018, Tarbouriech and Lazaric, 2019])
- Partially observable MDPs and beyond (e.g. [Jiang et al., 2017, Azizzadenesheli et al., 2016])

- 1 Setting the Stage
- 2 Lower Bounds
- 3 Optimism in Face of Uncertainty
- 4 Posterior Sampling
- 5 Asymptotically Optimal Algorithms
- 6 Extensions and Other Settings
- 7 Conclusion

Summary

Alg.	Asymptotic (ergodic)	Finite-time (comm.)
Lower bound	$\frac{C^2 S A}{\min_{s,a} \delta_{M^*}^*(s,a)} \ln(T)$	$\sqrt{D S A T}$
UCRL2B	$\frac{D^2 S^2 A}{\delta_g^*} \ln(T)$	$\sqrt{D S \Gamma A T \ln(T)}$
SCAL	$\frac{C^2 S^2 A}{\delta_g^*} \ln(T)$	$\sqrt{C S A T \ln(T)}$
TSDE	?	$C S \sqrt{A T \ln(T)}$
BKIA/DEL	$\frac{C^2 S A}{\min_{s,a} \delta_{M^*}^*(s,a)} \ln(T)$?

$$\blacksquare \Gamma = \max_{s,a} |\text{supp}(p(\cdot|s,a))|$$

$$\blacksquare D_M = \max_{s,s' \in \mathcal{S}} \min_{\pi: \mathcal{S} \rightarrow \mathcal{A}} \mathbb{E}[T_\pi^M(s,s')]$$

$$\blacksquare C \geq \text{sp}(h^*)$$


$$\blacksquare \delta_M^*(s,a) = L_M^* h_M^*(s) - L_M^a h_M^*(s)$$

$$\blacksquare \delta_g^* := g_M^* - \max_{s \in \mathcal{S}, \pi} \left\{ g_{M^*}^\pi(s) < g_M^* \right\}$$

Open Question: Summary

Alg.	Asymptotic (ergodic)	Finite-time (comm.)
Lower bound	$\frac{C^2 S A}{\min_{s,a} \delta_{M^*}^*(s,a)} \ln(T)$	$\sqrt{D S A T}$
UCRL2B	$\frac{D^2 S^2 A}{\delta_g^*} \ln(T)$	$\sqrt{D S T A T \ln(T)}$
SCAL	$\frac{C^2 S^2 A}{\delta_g^*} \ln(T)$	$\sqrt{C S T A T \ln(T)}$
TSDE	?	$C S \sqrt{A T \ln(T)}$ (Bayes)
BKIA/DEL	$\frac{C^2 S A}{\min_{s,a} \delta_{M^*}^*(s,a)} \ln(T)$?

Closing the gap between upper and lower bounds and settings (ergodic/asymptotic vs communicating/worst-case)

 Many lessons learned from bandit but need to deal with dynamical nature of the problem.

Open Questions

- Unifying finite-horizon, infinite-horizon regret and discounted PAC-MDP guarantees (e.g. [Dann et al., 2017])
- Model-based vs model-free (e.g. [Jin et al., 2018, Szepesvari et al., 2019])
- Scalable exp-exp (e.g. Bellemare et al. [2016], Tang and Agrawal [2018], Fortunato et al. [2017])

Open Questions: Model-free vs Model-based

Model-based exploration

- 👍 sample efficient (regret $O(\sqrt{T})$)
- 👎 solves an MDP at each episode ($O(S^2A)$)
- 👎 difficult to extend to function approximation

Model-free exploration

- 👎 sample inefficient (regret $O(T^{2/3})$?)
- 👍 simple update at each step ($O(1)$)
- 👍 easy to extend to function approximation

Open Questions: Model-free vs Model-based

Model-based exploration

- 👍 sample efficient (regret $O(\sqrt{T})$)
- 👎 solves an MDP at each episode ($O(S^2A)$)
- 👎 difficult to extend to function approximation

Model-free exploration

- 👎 sample inefficient (regret $O(T^{2/3})$?)
- 👍 simple update at each step ($O(1)$)
- 👍 easy to extend to function approximation

Sample and computationally efficient exploration algorithm? (see Jin et al. [2018])

Resources

Reinforcement Learning

■ Books

- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1994
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998
- Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. II*. Athena Scientific, 3rd edition, 2007
- Csaba Szepesvari. *Algorithms for Reinforcement Learning*. Morgan and Claypool Publishers, 2010

■ Courses (with good references for exploration)

- Nan Jiang. Cs598 statistical reinforcement learning.
<http://nanjiang.cs.illinois.edu/cs598/>
- Emma Brunskill. Cs234 reinforcement learning winter 2019.
<http://web.stanford.edu/class/cs234/index.html>
- Alessandro Lazaric. Mva reinforcement learning.
<http://chercheurs.lille.inria.fr/~lazaric/Webpage/Teaching.html>
- Alexandre Proutiere. Reinforcement learning: A graduate course.
http://www.it.uu.se/research/systems_and_control/education/2017/relearn/

Resources

Exploration-Exploitation and Regret Minimization

■ Books

- Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems.
Foundations and Trends® in Machine Learning, 5(1):1–122, 2012
- Tor Lattimore and Csaba Szepesvári. Bandit algorithms.
Pre-publication version, 2018.
URL <http://downloads.tor-lattimore.com/banditbook/book.pdf>



Thank you!

facebook

Artificial Intelligence Research



. \ |

- Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *COLT*, volume 19 of *JMLR Proceedings*, pages 1–26. JMLR.org, 2011.
- Yasin Abbasi-Yadkori and Csaba Szepesvári. Bayesian optimal control of smoothly parameterized systems. In *UAI*, pages 1–11. AUAI Press, 2015.
- Marc Abeille and Alessandro Lazaric. Thompson sampling for linear-quadratic control problems. In *AISTATS*, volume 54 of *Proceedings of Machine Learning Research*, pages 1246–1254. PMLR, 2017.
- Marc Abeille and Alessandro Lazaric. Improved regret bounds for thompson sampling in linear quadratic control problems. In *ICML*, volume 80 of *JMLR Workshop and Conference Proceedings*, pages 1–9. JMLR.org, 2018.
- Rajeev Agrawal. Adaptive control of markov chains under the weak accessibility. In *29th IEEE Conference on Decision and Control*, pages 1426–1431. IEEE, 1990.
- Shipra Agrawal and Randy Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. In *NIPS*, pages 1184–1194, 2017.
- Peter Auer and Ronald Ortner. Logarithmic online regret bounds for undiscounted reinforcement learning. In *NIPS*, pages 49–56. MIT Press, 2006.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 263–272. PMLR, 2017.
- Kamyar Azizzadenesheli, Alessandro Lazaric, and Animashree Anandkumar. Reinforcement learning of pomdps using spectral methods. In *COLT*, volume 49 of *JMLR Workshop and Conference Proceedings*, pages 193–256. JMLR.org, 2016.
- Peter L. Bartlett and Ambuj Tewari. REGAL: A regularization based algorithm for reinforcement learning in weakly communicating MDPs. In *UAI*, pages 35–42. AUAI Press, 2009.
- Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying count-based exploration and intrinsic motivation. In *NIPS*, pages 1471–1479, 2016.

- Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. II*. Athena Scientific, 3rd edition, 2007.
- Sergio Bittanti, Marco C Campi, et al. Adaptive control of linear time invariant systems: the “bet on the best” principle. *Communications in Information & Systems*, 6(4):299–320, 2006.
- Emma Brunskill. Cs234 reinforcement learning winter 2019. <http://web.stanford.edu/class/cs234/index.html>.
- Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- Apostolos N. Burnetas and Michael N. Katehakis. Optimal adaptive policies for markov decision processes. *Mathematics of Operations Research*, 22(1):222–255, 1997.
- Marco C Campi and PR Kumar. Adaptive linear quadratic gaussian control: the cost-biased approach revisited. *SIAM Journal on Control and Optimization*, 36(6):1890–1907, 1998.
- Sayak Ray Chowdhury and Aditya Gopalan. Online learning in kernelized markov decision processes. *CoRR*, abs/1805.08052, 2018.
- Alon Cohen, Avinatan Hassidim, Tomer Koren, Nevena Lazic, Yishay Mansour, and Kunal Talwar. Online linear quadratic control. In *ICML*, volume 80 of *JMLR Workshop and Conference Proceedings*, pages 1028–1037. JMLR.org, 2018.
- Alon Cohen, Tomer Koren, and Yishay Mansour. Learning Linear-Quadratic Regulators Efficiently with only $O(\sqrt{T})$ Regret. *arXiv e-prints*, art. arXiv:1902.06223, Feb 2019.
- Richard Combes, Stefan Magureanu, and Alexandre Proutière. Minimal exploration in structured stochastic bandits. In *NIPS*, pages 1761–1769, 2017.
- Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying PAC and regret: Uniform PAC bounds for episodic reinforcement learning. In *NIPS*, pages 5717–5727, 2017.
- Eyal Even-Dar, Sham. M. Kakade, and Yishay Mansour. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.

- Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Input perturbations for adaptive regulation and learning. *CoRR*, abs/1811.04258, 2018a.
- Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. On optimality of adaptive linear-quadratic regulators. *CoRR*, abs/1806.10749, 2018b.
- Sarah Filippi, Olivier Cappé, and Aurélien Garivier. Optimism in reinforcement learning and kullback-leibler divergence. *2010 48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 115–122, 2010.
- Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Rémi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy networks for exploration. *CoRR*, abs/1706.10295, 2017.
- Ronan Fruit and Alessandro Lazaric. Exploration-exploitation in mdps with options. In *AISTATS*, volume 54 of *Proceedings of Machine Learning Research*, pages 576–584. PMLR, 2017.
- Ronan Fruit, Matteo Pirota, Alessandro Lazaric, and Emma Brunskill. Regret minimization in mdps with options without prior knowledge. In *NIPS*, pages 3169–3179, 2017.
- Ronan Fruit, Matteo Pirota, and Alessandro Lazaric. Near optimal exploration-exploitation in non-communicating markov decision processes. In *NeurIPS*, pages 2998–3008, 2018a.
- Ronan Fruit, Matteo Pirota, Alessandro Lazaric, and Ronald Ortner. Efficient bias-span-constrained exploration-exploitation in reinforcement learning. In *ICML*, *Proceedings of Machine Learning Research*. PMLR, 2018b.
- Aurélien Garivier and Olivier Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *COLT*, volume 19 of *JMLR Proceedings*, pages 359–376. JMLR.org, 2011.
- Aditya Gopalan and Shie Mannor. Thompson sampling for learning parameterized markov decision processes. In *COLT*, volume 40 of *JMLR Workshop and Conference Proceedings*, pages 861–898. JMLR.org, 2015.

- Todd L Graves and Tze Leung Lai. Asymptotically efficient adaptive choice of control laws in controlled markov chains. *SIAM journal on control and optimization*, 35(3):715–743, 1997.
- Elad Hazan, Sham M. Kakade, Karan Singh, and Abby Van Soest. Provably Efficient Maximum Entropy Exploration. *arXiv e-prints*, art. arXiv:1812.02690, Dec 2018.
- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963. URL <http://www.jstor.org/stable/2282952>.
- Morteza Ibrahimi, Adel Javanmard, and Benjamin Van Roy. Efficient reinforcement learning for high dimensional linear quadratic systems. In *NIPS*, pages 2645–2653, 2012.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- Nan Jiang. Cs598 statistical reinforcement learning. <http://nanjiang.cs.illinois.edu/cs598/>.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E. Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 1704–1713. PMLR, 2017.
- Chi Jin, Zeyuan Allen-Zhu, Sébastien Bubeck, and Michael I. Jordan. Is q-learning provably efficient? In *NeurIPS*, pages 4868–4878, 2018.
- Sham Kakade, Mengdi Wang, and Lin F. Yang. Variance reduction methods for sublinear reinforcement learning. *CoRR*, abs/1802.09184, 2018.
- T.L Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4 – 22, 1985. ISSN 0196-8858. doi: [https://doi.org/10.1016/0196-8858\(85\)90002-8](https://doi.org/10.1016/0196-8858(85)90002-8). URL <http://www.sciencedirect.com/science/article/pii/0196885885900028>.
- K. Lakshmanan, Ronald Ortner, and Daniil Ryabko. Improved regret bounds for undiscounted continuous reinforcement learning. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 524–532. JMLR.org, 2015.

- Tor Lattimore and Csaba Szepesvári. Bandit algorithms. Pre-publication version, 2018. URL <http://downloads.tor-lattimore.com/banditbook/book.pdf>.
- Alessandro Lazaric. Mva reinforcement learning. <http://chercheurs.lille.inria.fr/~lazaric/Webpage/Teaching.html>.
- Shiau Hong Lim and Peter Auer. Autonomous exploration for navigating in mdps. In *COLT*, volume 23 of *JMLR Proceedings*, pages 40.1–40.24. JMLR.org, 2012.
- Gergely Neu, András György, Csaba Szepesvári, and András Antos. Online markov decision processes under bandit feedback. *IEEE Trans. Automat. Contr.*, 59(3):676–691, 2014.
- Jungseul Ok, Alexandre Proutière, and Damianos Tranos. Exploration in structured reinforcement learning. In *NeurIPS*, pages 8888–8896. 2018.
- Ronald Ortner. Online regret bounds for markov decision processes with deterministic transitions. *Theor. Comput. Sci.*, 411(29-30):2684–2695, 2010.
- Ronald Ortner and Daniil Ryabko. Online regret bounds for undiscounted continuous reinforcement learning. In *NIPS*, pages 1772–1780, 2012.
- Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. In *NIPS*, pages 1466–1474, 2014.
- Ian Osband and Benjamin Van Roy. On lower bounds for regret in reinforcement learning. *CoRR*, abs/1608.02732, 2016a.
- Ian Osband and Benjamin Van Roy. Posterior sampling for reinforcement learning without episodes. *CoRR*, abs/1608.02731, 2016b.
- Ian Osband and Benjamin Van Roy. Why is posterior sampling better than optimism for reinforcement learning? In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 2701–2710. PMLR, 2017.

- Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via posterior sampling. In *NIPS*, pages 3003–3011, 2013.
- Yi Ouyang, Mukul Gagrani, and Rahul Jain. Learning-based control of unknown linear systems with thompson sampling. *CoRR*, abs/1709.04047, 2017a. URL <http://arxiv.org/abs/1709.04047>.
- Yi Ouyang, Mukul Gagrani, Ashutosh Nayyar, and Rahul Jain. Learning unknown markov decision processes: A thompson sampling approach. In *NIPS*, pages 1333–1342, 2017b.
- Alexandre Proutiere. Reinforcement learning: A graduate course.
http://www.it.uu.se/research/systems_and_control/education/2017/relearn/.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1994.
- Jian Qian, Ronan Fruit, Matteo Pirota, and Alessandro Lazaric. Exploration bonus for regret minimization in undiscounted discrete and continuous markov decision processes. *CoRR*, 2018.
- Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- Malcolm Strens. A bayesian framework for reinforcement learning. In *In Proceedings of the Seventeenth International Conference on Machine Learning*, pages 943–950. ICML, 2000.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- Csaba Szepesvari. *Algorithms for Reinforcement Learning*. Morgan and Claypool Publishers, 2010.
- Csaba Szepesvari, Nevena Lazic, and Yasin Abbasi-Yadkori. Model-free linear quadratic control via reduction to expert prediction. In *AISTATS*, 2019.
- Mohammad Sadegh Talebi and Odalric-Ambrym Maillard. Variance-aware regret bounds for undiscounted reinforcement learning in mdps. In *ALT*, volume 83 of *Proceedings of Machine Learning Research*, pages 770–805. PMLR, 2018.

- Yunhao Tang and Shipra Agrawal. Exploration by distributional reinforcement learning. *CoRR*, abs/1805.01907, 2018.
- Jean Tarbouriech and Alessandro Lazaric. Active Exploration in Markov Decision Processes. *arXiv e-prints*, art. arXiv:1902.11199, Feb 2019.
- Ambuj Tewari and Peter L. Bartlett. Optimistic linear programming gives logarithmic regret for irreducible mdps. In *NIPS*, pages 1505–1512. Curran Associates, Inc., 2007.
- Georgios Theocharous, Zheng Wen, Yasin Abbasi, and Nikos Vlassis. Scalar posterior sampling with applications. In *NeurIPS*, pages 7696–7704, 2018.
- William R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- Tsachy Weissman, Erik Ordentlich, Gadiel Seroussi, Sergio Verdú, and Marcelo J. Weinberger. Inequalities for the L1 deviation of the empirical distribution. 2003.
- Andrea Zanette and Emma Brunskill. Problem dependent reinforcement learning bounds which can identify bandit structure in mdps. In *ICML*, volume 80 of *JMLR Workshop and Conference Proceedings*, pages 5732–5740. JMLR.org, 2018.